

## Chapter 1

# DESCENT METHODS FOR NONNEGATIVE MATRIX FACTORIZATION

Ngoc-Diep Ho, Paul Van Dooren and Vincent D. Blondel

*CESAME, Université catholique de Louvain,  
Av. Georges Lemaitre 4, B-1348 Louvain-la-Neuve, Belgium.  
Tel : +32(10) 47 22 64 Fax : +32(10) 47 21 80  
{ngoc.ho, paul.vandooren, vincent.blondel}@uclouvain.be*

**Abstract** In this paper, we present several descent methods that can be applied to nonnegative matrix factorization and we analyze a recently developed fast block coordinate method called Rank-one Residue Iteration (RRI). We also give a comparison of these different methods and show that the new block coordinate method has better properties in terms of approximation error and complexity. By interpreting this method as a rank-one approximation of the residue matrix, we prove that it *converges* and also extend it to the nonnegative tensor factorization and introduce some variants of the method by imposing some additional controllable constraints such as: sparsity, discreteness and smoothness.

**Keywords:** Algorithm, Nonnegative matrix, Factorization

## 1. Introduction

Linear algebra has become a key tool in almost all modern techniques for data analysis. Most of these techniques make use of linear subspaces represented by eigenvectors of a particular matrix. In this paper, we consider a set of  $n$  data points  $a_1, a_2, \dots, a_n$ , where each point is a real vector of size  $m$ ,  $a_i \in \mathbb{R}^m$ . We then approximate these data points by linear combinations of  $r$  basis vectors  $u_i \in \mathbb{R}^m$ :

$$a_i \approx \sum_{j=1}^r v_{ij} u_j, \quad v_{ij} \in \mathbb{R}, \quad u_j \in \mathbb{R}^m.$$

This can be rewritten in matrix form as  $A \approx UV^T$ , where  $a_i$  and  $u_i$  are respectively the columns of  $A$  and  $U$  and the  $v_{ij}$ 's are the elements of  $V$ . Optimal solutions of this approximation in terms of the Euclidean (or Frobenius) norm can be obtained by the Singular Value Decomposition (SVD) [13].

In many cases, data points are constrained to a subset of  $\mathbb{R}^m$ . For example, light intensities, concentrations of substances, absolute temperatures are, by their nature, nonnegative (or even positive) and lie in the nonnegative orthant  $\mathbb{R}_+^m$ . The input matrix  $A$  then becomes elementwise nonnegative and it is then natural to constrain the basis vectors  $v_i$  and the coefficients  $v_{ij}$  to be nonnegative as well. In order to satisfy this constraint, we need to approximate the columns of  $A$  by the following additive model:

$$a_i \approx \sum_{j=1}^r v_{ij} u_j, \quad v_{ij} \in \mathbb{R}_+, \quad u_j \in \mathbb{R}_+^m.$$

where the  $v_{ij}$  coefficients and  $u_j$  vectors are nonnegative,  $v_{ij} \in \mathbb{R}_+$ ,  $u_j \in \mathbb{R}_+^m$ .

Many algorithms have been proposed to find such a representation, which is referred to as a Nonnegative Matrix Factorization (NMF). The earliest algorithms were introduced by Paatero [25, 26]. But the topic became quite popular with the publication of the algorithm of Lee and Seung in 1999 [20] where multiplicative rules were introduced to solve the problem. This algorithm is very simple and elegant but it lacks a complete convergence analysis. Other methods and variants can be found in [23], [21], [17].

The quality of the approximation is often measured by a distance. Two popular choices are the Euclidean (Frobenius) norm and the generalized Kullback-Leibler divergence. In this paper, we focus on the Euclidean distance and we investigate descent methods for this measure. One characteristic of descent methods is their monotonic decrease until they reach a stationary point. This point maybe located in the interior of the nonnegative orthant or on its boundary. In the second case, the constraints become active and may prohibit any further decrease of the distance measure. This is a key issue to be analyzed for any descent method.

In this paper,  $\mathbb{R}_+^m$  denotes the set of nonnegative real vectors (elementwise) and  $[v]_+$  the projection of the vector  $v$  on  $\mathbb{R}_+^m$ . We use  $v \geq 0$  and  $A \geq 0$  to denote nonnegative vectors and matrices and  $v > 0$  and  $A > 0$  to denote positive vectors and matrices.  $A \circ B$  and  $\frac{[A]}{[B]}$  are respectively

the Hadamard (elementwise) product and quotient.  $A_{:,i}$  and  $A_{i,:}$  are the  $i^{th}$  column and  $i^{th}$  row of  $A$ .

This paper is an extension of the internal report [15], where we proposed to decouple the problem based on rank one approximations to create a new algorithm called Rank-one Residue Iteration (RRI). During the revision of this report, we were informed that essentially the same algorithm was independently proposed and published in [8] under the name Hierarchical Alternative Least Squares (HALS). But the present paper gives several additional results wherein the major contributions are *the convergence proof* of the method and its *extensions* to many practical situations and constraints. The paper also compares a selection of some recent descent methods from the literature and aims at providing a survey of such methods for nonnegative matrix factorizations. For that reason, we try to be self-contained and hence recall some well-known results. We also provide short proofs when useful for a better understanding of the rest of the paper.

We first give a short introduction of low rank approximations, both unconstrained and constrained. In Section 3 we discuss error bounds of various approximations and in Section 4 we give a number of descent methods for Nonnegative Matrix Factorizations. In Section 5 we describe the method based on successive rank one approximations. This method is then also extended to approximate higher order tensor and to take into account other constraints than nonnegativity. In Section 5.0 we discuss various regularization methods and in Section 6, we present numerical experiments comparing the different methods. We end with some concluding remarks.

## 2. Low-rank matrix approximation

Low-rank approximation is a special case of matrix nearness problem [14]. When only a rank constraint is imposed, the optimal approximation with respect to the Frobenius norm can be obtained from the Singular Value Decomposition.

We first investigate the problem without the nonnegativity constraint on the low-rank approximation. This is useful for understanding properties of the approximation when the nonnegativity constraints are imposed but inactive. We begin with the well-known Eckart-Young Theorem.

**Theorem 2.1** (Eckart-Young). *Let  $A \in \mathbb{R}^{m \times n}$  ( $m \geq n$ ) have the singular value decomposition*

$$A = P\Sigma Q^T, \quad \Sigma = \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix}$$

where  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$  are the singular values of  $A$  and where  $P \in \mathbb{R}^{m \times m}$  and  $Q \in \mathbb{R}^{n \times n}$  are orthogonal matrices. Then for  $1 \leq r \leq n$ , the matrix

$$A_r = P\Sigma_r Q^T, \quad \Sigma_r = \begin{pmatrix} \sigma_1 & 0 & \dots & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & & \vdots \\ 0 & 0 & \dots & \sigma_r & \dots & 0 \\ \vdots & \vdots & & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \dots & 0 \end{pmatrix}$$

is a global minimizer of the problem

$$\min_{B \in \mathbb{R}^{m \times n} \text{ rank}(B) \leq r} \frac{1}{2} \|A - B\|_F^2 \quad (1.1)$$

and its error is

$$\frac{1}{2} \|A - B\|_F^2 = \frac{1}{2} \sum_{i=r+1}^n \sigma_i^2.$$

Moreover, if  $\sigma_r > \sigma_{r+1}$  then  $A_r$  is the unique global minimizer.

The proof and other implications can be found for instance in [13]. The columns of  $P$  and  $Q$  are called singular vectors of  $A$ , in which vectors corresponding to the largest singular values are referred to as the dominant singular vectors.

Let us now look at the following modified problem

$$\min_{X \in \mathbb{R}^{m \times r} \ Y \in \mathbb{R}^{n \times r}} \frac{1}{2} \|A - XY^T\|_F^2, \quad (1.2)$$

where the rank constraint is implicit in the product  $XY^T$  since the dimensions of  $X$  and  $Y$  guarantee that  $\text{rank}(XY^T) \leq r$ . Conversely, every matrix of rank less than  $r$  can be trivially rewritten as a product

$XY^T$ , where  $X \in \mathbb{R}^{m \times r}$  and  $Y \in \mathbb{R}^{n \times r}$ . Therefore Problems (1.1) and (1.2) are equivalent. But even when the product  $A_r = XY^T$  is unique, the pairs  $(XR^T, YR^{-1})$  with  $R$  invertible, yield the same product  $XY^T$ . In order to avoid this, we can always choose  $X$  and  $Y$  such that

$$X = PD^{\frac{1}{2}} \text{ and } Y = QD^{\frac{1}{2}}, \quad (1.3)$$

where  $P^T P = I_{r \times r}$ ,  $Q^T Q = I_{r \times r}$  and  $D$  is  $r \times r$  nonnegative diagonal matrix. Doing this is equivalent to computing a compact SVD decomposition of the product  $A_r = XY^T = PDQ^T$ .

As usual for optimization problems, we calculate the gradient with respect to  $X$  and  $Y$  and set them equal to 0.

$$\nabla_X = XY^T Y - AY = 0 \quad \nabla_Y = YX^T X - A^T X = 0. \quad (1.4)$$

If we then premultiply  $A^T$  with  $\nabla_X$  and  $A$  with  $\nabla_Y$ , we obtain

$$(A^T A)Y = (A^T X)Y^T Y \quad (AA^T)X = (AY)X^T X. \quad (1.5)$$

Replacing  $A^T X = YX^T X$  and  $AY = XY^T Y$  into (1.5) yields

$$(A^T A)Y = YX^T XY^T Y \quad (AA^T)X = XY^T YX^T X. \quad (1.6)$$

Replacing (1.3) into (1.6) yields

$$(A^T A)QD^{\frac{1}{2}} = QDP^T PDQ^T QD^{\frac{1}{2}} \text{ and } (AA^T)PD^{\frac{1}{2}} = PDQ^T QDP^T PD^{\frac{1}{2}}.$$

When  $D$  is invertible, this finally yields

$$(A^T A)Q = QD^2 \text{ and } (AA^T)P = PD^2.$$

This shows that the columns of  $P$  and  $Q$  are singular vectors and  $D_{ii}$ 's are nonzero singular values of  $A$ . Notice that if  $D$  is singular, one can throw away the corresponding columns of  $P$  and  $Q$  and reduce it to a smaller-rank approximation with the same properties. Without loss of generality, we therefore can focus on approximations of Problem (1.2) which are of exact rank  $r$ . We can summarize the above reasoning in the following theorem.

**Theorem 2.2.** *Let  $A \in \mathbb{R}^{m \times n}$  ( $m > n$  and  $\text{rank}(A) = t$ ). If  $A_r$  ( $1 \leq r \leq t$ ) is a rank  $r$  stationary point of Problem 1.2, then there exists two orthogonal matrices  $P \in \mathbb{R}^{m \times m}$  and  $Q \in \mathbb{R}^{n \times n}$  such that:*

$$A = P\hat{\Sigma}Q^T \text{ and } A_r = P\hat{\Sigma}_rQ^T$$

where

$$\hat{\Sigma} = \begin{pmatrix} \hat{\sigma}_1 & 0 & \dots & 0 \\ 0 & \hat{\sigma}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \hat{\sigma}_n \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix}, \quad \hat{\Sigma}_r = \begin{pmatrix} \hat{\sigma}_1 & 0 & \dots & 0 & \dots & 0 \\ 0 & \hat{\sigma}_2 & \dots & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & & \vdots \\ 0 & 0 & \dots & \hat{\sigma}_r & \dots & 0 \\ \vdots & \vdots & & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \dots & 0 \end{pmatrix}$$

and the  $\hat{\sigma}_i$ 's are unsorted singular values of  $A$ . Moreover, the approximation error is:

$$\frac{1}{2}\|A - A_r\|_F^2 = \frac{1}{2} \sum_{i=r+1}^t \hat{\sigma}_i^2.$$

This result shows that, if the singular values are all different, there are  $\frac{n!}{r!(n-r)!}$  possible stationary points  $A_r$ . When there are multiple singular values, there will be infinitely many stationary points  $A_r$  since there are infinitely many singular subspaces. The next result will identify the minima among all stationary points. Other stationary points are saddle points whose every neighborhood contains both smaller and higher points.

**Theorem 2.3.** *The only minima of Problem 1.2 are given by Theorem 2.1 and are global minima. All other stationary points are saddle points.*

*Proof.* Let us assume that  $A_r$  is a stationary point given by Theorem 2.2 but not by Theorem 2.1. Then there always exists a permutation of the columns of  $P$  and  $Q$ , and of the diagonal elements of  $\hat{\Sigma}$  and  $\hat{\Sigma}_r$  such that  $\hat{\sigma}_{r+1} > \hat{\sigma}_r$ . We then construct two points in the  $\epsilon$ -neighborhood of  $A_r$  that yield an increase and a decrease, respectively, of the distance measure. They are obtained by taking:

$$\bar{\Sigma}_r(\epsilon) = \begin{pmatrix} \hat{\sigma}_1 + \epsilon & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \dots & \vdots \\ 0 & \dots & \hat{\sigma}_r & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{pmatrix}, \quad \bar{A}_r(\epsilon) = P\bar{\Sigma}_r(\epsilon)Q^T$$

and

$$\underline{\Sigma}_r(\epsilon) = \begin{pmatrix} \hat{\sigma}_1 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \dots & \vdots \\ 0 & \dots & \hat{\sigma}_r & \epsilon\sqrt{\hat{\sigma}_r} & \vdots & 0 \\ 0 & \dots & \epsilon\sqrt{\hat{\sigma}_r} & \epsilon^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \dots & 0 \end{pmatrix}, \quad \underline{A}_r(\epsilon) = P\underline{\Sigma}_r(\epsilon)Q^T.$$

Clearly  $\overline{A}_r(\epsilon)$  and  $\underline{A}_r(\epsilon)$  are of rank  $r$ . Evaluating the distance measure yields

$$\begin{aligned} \|A - \underline{A}_r(\epsilon)\|_F^2 &= 2\hat{\sigma}_r\epsilon^2 + (\hat{\sigma}_{r+1} - \epsilon^2)^2 + \sum_{i=r+2}^t \hat{\sigma}_i^2 \\ &= \epsilon^2[\epsilon^2 - 2(\hat{\sigma}_{r+1} - \hat{\sigma}_r)] + \sum_{i=r+1}^t \hat{\sigma}_i^2 \\ &< \sum_{i=r+1}^t \hat{\sigma}_i^2 = \|A - A_r\|_F^2 \end{aligned}$$

for all  $\epsilon \in (0, \sqrt{2(\hat{\sigma}_{r+1} - \hat{\sigma}_r)})$  and

$$\|A - \overline{A}_r(\epsilon)\|_F^2 = \epsilon^2 + \sum_{i=r+1}^t \hat{\sigma}_i^2 > \sum_{i=r+1}^t \hat{\sigma}_i^2 = \|A - A_r\|_F^2$$

for all  $\epsilon > 0$ . Hence, for an arbitrarily small positive  $\epsilon$ , we obtain

$$\|A - \underline{A}_r(\epsilon)\|_F^2 < \|A - A_r\|_F^2 < \|A - \overline{A}_r(\epsilon)\|_F^2$$

which shows that  $A_r$  is a saddle point of the distance measure.  $\square$

When we add a nonnegativity constraint in the next section, the results of this section will help to identify stationary points at which all the nonnegativity constraints are inactive.

### 3. Nonnegativity constraint

In this section, we investigate the problem of Nonnegative Matrix Factorization. This problem differs Problem 1.2 in the previous section because of the additional nonnegativity constraints on the factors. We first discuss the effects of adding such a constraint. By doing so, the

problem is no longer easy because of the existence of local minima at the boundary of the nonnegative orthant. Determining the lowest minimum among these minima is far from trivial. On the other hand, a minimum that coincides with a minimum of the unconstrained problem (i.e. Problem 1.2) may be easily reached by standard descent methods, as we will see.

**Problem 1** (Nonnegative matrix factorization - NMF). *Given a  $m \times n$  nonnegative matrix  $A$  and an integer  $r < \min(m, n)$ , solve*

$$\min_{U \in \mathbb{R}_+^{m \times r} \quad V \in \mathbb{R}_+^{n \times r}} \frac{1}{2} \|A - UV^T\|_F^2.$$

Where  $r$  is called the reduced rank. From now on,  $m$  and  $n$  will be used to denote the size of the target matrix  $A$  and  $r$  is the reduced rank of a factorization.

We rewrite the nonnegative matrix factorization as a standard nonlinear optimization problem:

$$\min_{-U \leq 0 \quad -V \leq 0} \frac{1}{2} \|A - UV^T\|_F^2.$$

The associated Lagrangian function is

$$L(U, V, \mu, \nu) = \frac{1}{2} \|A - UV^T\|_F^2 - \mu \circ U - \nu \circ V,$$

where  $\mu$  and  $\nu$  are two *matrices* of the same size of  $U$  and  $V$ , respectively, containing the Lagrange multipliers associated with the nonnegativity constraints  $U_{ij} \geq 0$  and  $V_{ij} \geq 0$ . Then the Karush-Kuhn-Tucker conditions for the nonnegative matrix factorization problem say that if  $(U, V)$  is a local minimum, then there exist  $\mu_{ij} \geq 0$  and  $\nu_{ij} \geq 0$  such that:

$$U \geq 0 \quad , \quad V \geq 0, \tag{1.7}$$

$$\nabla L_U = 0 \quad , \quad \nabla L_V = 0, \tag{1.8}$$

$$\mu \circ U = 0 \quad , \quad \nu \circ V = 0. \tag{1.9}$$

Developing (1.8) we have:

$$AV - UV^T V - \mu = 0, \quad A^T U - V U^T U - \nu = 0$$

or

$$\mu = -(UV^T V - AV), \quad \nu = -(V U^T U - A^T U).$$



Combining this with  $\mu_{ij} \geq 0$ ,  $\nu_{ij} \geq 0$  and (1.9) gives the following conditions:

$$U \geq 0 \quad , \quad V \geq 0, \quad (1.10)$$

$$\nabla F_U = UV^T V - AV \geq 0 \quad , \quad \nabla F_V = VU^T U - A^T U \geq 0, \quad (1.11)$$

$$U \circ (UV^T V - AV) = 0 \quad , \quad V \circ (VU^T U - A^T U) = 0, \quad (1.12)$$

where the corresponding Lagrange multipliers for  $U$  and  $V$  are also the gradient of  $F$  with respect to  $U$  and  $V$ . Since the Euclidean distance is not convex with respect to both variables  $U$  and  $V$  at the same time, these conditions are only necessary. This is implied because of the existence of saddle points and maxima. We then call all the points that satisfy the above conditions, the stationary points.

**Definition 1** (NMF stationary point). *We call  $(U, V)$  a stationary point of the NMF Problem if and only if  $U$  and  $V$  satisfy the KKT conditions (1.10), (1.11) and (1.12).*

Alternatively, a stationary point  $(U, V)$  of the NMF problem can also be defined by using the following necessary condition (see for example [4]) on the convex sets  $\mathbb{R}_+^{m \times r}$  and  $\mathbb{R}_+^{n \times r}$ , that is

$$\left\langle \begin{pmatrix} \nabla F_U \\ \nabla F_V \end{pmatrix}, \begin{pmatrix} X - U \\ Y - V \end{pmatrix} \right\rangle \geq 0, \quad \forall X \in \mathbb{R}_+^{m \times r}, Y \in \mathbb{R}_+^{n \times r}, \quad (1.13)$$

which can be shown to be equivalent to the KKT conditions (1.10), (1.11) and (1.12). Indeed, it is trivial that the KKT conditions imply (1.13). And by carefully choosing different values of  $X$  and  $Y$  from (1.13), one can easily prove that the KKT conditions hold.

There are two values of reduced rank  $r$  for which we can trivially identify the global solution which are  $r = 1$  and  $r = \min(m, n)$ . For  $r = 1$ , a pair of dominant singular vectors are a global minimizer. And for  $r = \min(m, n)$ ,  $(U = A, V = I)$  is a global minimizer. Since most of existing methods for the nonnegative matrix factorization are descent algorithms, we should pay attention to all local minimizers. For the rank-one case, they can easily be characterized.

## Rank one case

The rank-one NMF problem of a nonnegative matrix  $A$  can be rewritten as

$$\min_{u \in \mathbb{R}_+^m, v \in \mathbb{R}_+^n} \frac{1}{2} \|A - uv^T\|_F^2 \quad (1.14)$$

and a complete analysis can be carried out. It is well known that any pair of nonnegative Perron vectors of  $AA^T$  and  $A^T A$  yields a global minimizer

of this problem, but we can also show that the *only* stationary points of (1.14) are given by such vectors. The following theorem excludes the case where  $u = 0$  and/or  $v = 0$ .

**Theorem 3.1.** *The pair  $(u, v)$  is a local minimizer of (1.14) if and only if  $u$  and  $v$  are nonnegative eigenvectors of  $AA^T$  and  $A^T A$  respectively of the eigenvalue  $\sigma = \|u\|_2^2 \|v\|_2^2$ .*

*Proof.* The *if part* easily follows from Theorem 2.2. For the *only if part* we proceed as follows. Without loss of generality, we can permute the rows and columns of  $A$  such that the corresponding vectors  $u$  and  $v$  are partitioned as  $(u_+ \ 0)^T$  and  $(v_+ \ 0)^T$  respectively, where  $u_+, v_+ > 0$ . Partition the corresponding matrix  $A$  conformably as follows

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix},$$

then from (1.11) we have

$$\begin{pmatrix} u_+ v_+^T & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} v_+ \\ 0 \end{pmatrix} - \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \begin{pmatrix} v_+ \\ 0 \end{pmatrix} \geq 0$$

and

$$\begin{pmatrix} v_+ u_+^T & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} u_+ \\ 0 \end{pmatrix} - \begin{pmatrix} A_{11}^T & A_{21}^T \\ A_{12}^T & A_{22}^T \end{pmatrix} \begin{pmatrix} u_+ \\ 0 \end{pmatrix} \geq 0$$

implying that  $A_{21}v_+ \leq 0$  and  $A_{12}^T u_+ \leq 0$ . Since  $A_{21}, A_{12} \geq 0$  and  $u_+, v_+ > 0$ , we can conclude that  $A_{12} = 0$  and  $A_{21} = 0$ . Then from (1.12) we have:

$$u_+ \circ (\|v_+\|_2^2 u_+ - A_{11}v_+) = 0 \text{ and } v_+ \circ (\|u_+\|_2^2 v_+ - A_{11}^T u_+) = 0.$$

Since  $u_+, v_+ > 0$ , we have:

$$\|v_+\|_2^2 u_+ = A_{11}v_+ \text{ and } \|u_+\|_2^2 v_+ = A_{11}^T u_+$$

or

$$\|u_+\|_2^2 \|v_+\|_2^2 u_+ = A_{11}A_{11}^T u_+ \text{ and } \|u_+\|_2^2 \|v_+\|_2^2 v_+ = A_{11}^T A_{11} v_+.$$

Setting  $\sigma = \|u_+\|_2^2 \|v_+\|_2^2$  and using the block *diagonal* structure of  $A$  yields the desired result.  $\square$

Theorem 3.1 guarantees that all stationary points of the rank-one case are nonnegative singular vectors of a submatrix of  $A$ . These results imply that a global minimizer of the rank-one NMF can be calculated

correctly based on the largest singular value and corresponding singular vectors of the matrix  $A$ .

For ranks other than 1 and  $\min(m, n)$ , there are no longer trivial stationary points. In the next section, we try to derive some simple characteristics of the local minima of the nonnegative matrix factorization.

The KKT conditions (1.12) help to characterize the stationary points of the NMF problem. Summing up all the elements of one of the conditions (1.12), we get:

$$\begin{aligned} 0 &= \sum_{ij} (U \circ (UV^T V - AV))_{ij} \\ &= \langle U, UV^T V - AV \rangle \\ &= \langle UV^T, UV^T - A \rangle. \end{aligned} \tag{1.15}$$

From that, we have some simple characteristics of the NMF solutions:

**Theorem 3.2.** *Let  $(U, V)$  be a stationary point of the NMF problem, then  $UV^T \in \mathcal{B}(\frac{A}{2}, \frac{1}{2}\|A\|_F)$ , the ball centered at  $\frac{A}{2}$  and with radius  $= \frac{1}{2}\|A\|_F$ .*

*Proof.* From (1.15) it immediately follows that

$$\left\langle \frac{A}{2} - UV^T, \frac{A}{2} - UV^T \right\rangle = \left\langle \frac{A}{2}, \frac{A}{2} \right\rangle$$

which implies

$$UV^T \in \mathcal{B}\left(\frac{A}{2}, \frac{1}{2}\|A\|_F\right).$$

**Theorem 3.3.** *Let  $(U, V)$  be a stationary of the NMF problem, then*

$$\frac{1}{2}\|A - UV^T\|_F^2 = \frac{1}{2}(\|A\|_F^2 - \|UV^T\|_F^2).$$

*Proof.* From (1.15), we have  $\langle UV^T, A \rangle = \langle UV^T, UV^T \rangle$ . Therefore,

$$\begin{aligned} \frac{1}{2}\langle A - UV^T, A - UV^T \rangle &= \frac{1}{2}(\|A\|_F^2 - 2\langle UV^T, A \rangle + \|UV^T\|_F^2) \\ &= \frac{1}{2}(\|A\|_F^2 - \|UV^T\|_F^2). \end{aligned}$$

□

Theorem 3.3 also suggests that at a stationary point  $(U, V)$  of the NMF problem, we should have  $\|A\|_F^2 \geq \|UV^T\|_F^2$ . This norm inequality

can be also found in [7] for less general cases where we have  $\nabla F_U = 0$  and  $\nabla F_V = 0$  at a stationary point. For this particular class of NMF stationary point, all the nonnegativity constraints on  $U$  and  $V$  are inactive. And all such stationary points are also stationary points of the unconstrained problem, characterized by Theorem 2.2.

We have seen in Theorem 2.2 that, for the unconstrained least-square problem the only stable stationary points are in fact global minima. Therefore, if the stationary points of the constrained problem are inside the nonnegative orthant (i.e. all constraints are inactive), we can then probably reach the global minimum of the NMF problem. This can be expected because the constraints may no longer prohibit the descent of the update.

Let  $A_r$  be the optimal rank- $r$  approximation of a nonnegative matrix  $A$ , which we obtain from the singular value decomposition, as indicated in Theorem 2.2. Then we can easily construct its nonnegative part  $[A_r]_+$ , which is obtained from  $A_r$  by just setting all its negative elements equal to zero. This is in fact the closest matrix in the cone of nonnegative matrices to the matrix  $A_r$ , in the Frobenius norm (in that sense, it is its projection on that cone). We now derive some bounds for the error  $\|A - [A_r]_+\|_F$ .

**Theorem 3.4.** *Let  $A_r$  be the best rank  $r$  approximation of a nonnegative matrix  $A$ , and let  $[A_r]_+$  be its nonnegative part, then*

$$\|A - [A_r]_+\|_F \leq \|A - A_r\|_F.$$

*Proof.* This follows easily from the convexity of the cone of nonnegative matrices. Since both  $A$  and  $[A_r]_+$  are nonnegative and since  $[A_r]_+$  is the closest matrix in that cone to  $A_r$  we immediately obtain the inequality

$$\|A - A_r\|_F^2 \geq \|A - [A_r]_+\|_F^2 + \|A_r - [A_r]_+\|_F^2 \geq \|A - [A_r]_+\|_F^2$$

from which the result readily follows.  $\square$

The approximation  $[A_r]_+$  has the merit of requiring as much storage as a rank  $r$  approximation, even though its rank is larger than  $r$  whenever  $A_r \neq [A_r]_+$ . We will look at the quality of this approximation in Section 6. If we now compare this bound with the nonnegative approximations then we obtain the following inequalities. Let  $U_* V_*^T$  be an optimal nonnegative rank  $r$  approximation of  $A$  and let  $UV^T$  be any stationary point of the KKT conditions for a nonnegative rank  $r$  approximation, then we have :

$$\|A - [A_r]_+\|_F^2 \leq \|A - A_r\|_F^2 = \sum_{i=r+1}^n \sigma_i^2 \leq \|A - U_* V_*^T\|_F^2 \leq \|A - UV^T\|_F^2.$$

For more implications of the NMF problem, see [16].

#### 4. Existing descent algorithms

We focus on descent algorithms that guarantee a non increasing update at each iteration. Based on the search space, we have two categories: *Full-space search* and *(Block) Coordinate search*.

Algorithms in the former category try to find updates for both  $U$  and  $V$  at the same time. This requires a search for a descent direction in the  $(m+n)r$ -dimensional space. Note also that the NMF problem in this full space is not convex but the optimality conditions may be easier to achieve.

Algorithms in the latter category, on the other hand, find updates for each (block) coordinate in order to guarantee the descent of the objective function. Usually, search subspaces are chosen to make the objective function convex so that efficient methods can be applied. Such a simplification might lead to the loss of some convergence properties. Most of the algorithms use the following column partitioning:

$$\frac{1}{2}\|A - UV^T\|_F^2 = \frac{1}{2} \sum_{i=1}^n \|A_{:,i} - U(V_{i,:})^T\|_2^2, \quad (1.16)$$

which shows that one can minimize with respect to each of the rows of  $V$  independently. The problem thus decouples into smaller convex problems. This leads to the solution of quadratic problems of the form

$$\min_{v \geq 0} \frac{1}{2} \|a - Uv\|_2^2. \quad (1.17)$$

Updates for the rows of  $V$  are then alternated with updates for the rows of  $U$  in a similar manner by transposing  $A$  and  $UV^T$ .

Independent on the search space, most of algorithms use the Projected Gradient scheme for which three basic steps are carried out in each iteration:

- Calculating the gradient  $\nabla F(x^k)$ ,
- Choosing the step size  $\alpha^k$ ,
- Projecting the update on the nonnegative orthant

$$x^{k+1} = [x^k - \alpha^k \nabla F(x^k)]_+,$$

where  $x^k$  is the variable in the selected search space. The last two steps can be merged in one iterative process and must guarantee a sufficient decrease of the objective function as well as the nonnegativity of the new point.

## Multiplicative rules (Mult)

Multiplicative rules were introduced in [20]. The algorithm applies a block coordinate type search and uses the above column partition to formulate the updates. A special feature of this method is that the step size is calculated for each element of the vector. For the elementary problem (1.17) it is given by

$$v^{k+1} = v^k - \alpha^k \circ \nabla F(v^{k+1}) = v^k \circ \frac{[U^T a]}{[U^T U v^k]}$$

where  $[\alpha^k]_i = \frac{v_i}{[U^T U v]_i}$ . Applying this to all rows of  $V$  and  $U$  gives the updating rule of Algorithm 1 to compute

$$(U^*, V^*) = \underset{U \geq 0 \ V \geq 0}{\operatorname{argmin}} \|A - UV^T\|_F^2.$$

---

### Algorithm 1 (Mult)

---

- 1: Initialize  $U^0$ ,  $V^0$  and  $k = 0$
  - 2: **repeat**
  - 3:    $U^{k+1} = U^k \circ \frac{[AV^k]}{[U^k(V^k)^T(V^k)]}$
  - 4:    $V^{k+1} = V^k \circ \frac{[A^T U^{k+1}]}{[V^k(U^{k+1})^T(U^{k+1})]}$
  - 5:    $k = k + 1$
  - 6: **until** Stopping condition
- 

These updates guarantee automatically the nonnegativity of the factors but may fail to give a sufficient decrease of the objective function. It may also get stuck in a non-stationary point and hence suffer from a poor convergence. Variants can be found in [21, 24].

## Line search using Armijo criterion (Line)

In order to ensure a sufficient descent, the following projected gradient scheme with Armijo criterion [23, 22] can be applied to minimize

$$x^* = \underset{x}{\operatorname{argmin}} F(x).$$

Algorithm 2 needs two parameters  $\sigma$  and  $\beta$  that may affect its convergence. It requires only the gradient information, and is applied in [23] for two different strategies : for the whole space  $(U, V)$  (Algorithm FLine) and for  $U$  and  $V$  separately in an alternating fashion (Algorithm

---

**Algorithm 2** (Line)

---

```

1: Initialize  $x^0$ ,  $\sigma$ ,  $\beta$ ,  $\alpha_0 = 1$  and  $k = 1$ 
2: repeat
3:    $\alpha_k = \alpha_{k-1}$ 
4:    $y = [x^k - \alpha_k \nabla F(x^k)]_+$ 
5:   if  $F(y) - F(x^k) > \sigma \langle \nabla F(x^k), y - x^k \rangle$  then
6:     repeat
7:        $\alpha_k = \alpha_k \cdot \beta$ 
8:        $y = [x^k - \alpha_k \nabla F(x^k)]_+$ 
9:     until  $F(y) - F(x^k) \leq \sigma \langle \nabla F(x^k), y - x^k \rangle$ 
10:  else
11:    repeat
12:       $lasty = y$ 
13:       $\alpha_k = \alpha_k / \beta$ 
14:       $y = [x^k - \alpha_k \nabla F(x^k)]_+$ 
15:    until  $F(y) - F(x^k) > \sigma \langle \nabla F(x^k), y - x^k \rangle$ 
16:     $y = lasty$ 
17:  end if
18:   $x^{k+1} = y$ 
19:   $k = k + 1$ 
20: until Stopping condition

```

---

CLine). With a good choice of parameters ( $\sigma = 0.01$  and  $\beta = 0.1$ ) and a good strategy of alternating between variables, it was reported in [23] to be the faster than the multiplicative rules.

### Projected gradient with first-order approximation (FO)

In order to find the solution to

$$x^* = \operatorname{argmin}_x F(x)$$

we can also approximate at each iteration the function  $F(X)$  using:

$$\tilde{F}(x) = F(x^k) + \left\langle \nabla_x F(x^k), x - x^k \right\rangle + \frac{L}{2} \|x^k - x\|_2^2,$$

where  $L$  is a Lipschitz constant satisfying  $F(x) \leq \tilde{F}(x)$ ,  $\forall x$ . Because of this inequality, the solution of the following problem

$$x_{k+1} = \operatorname{argmin}_{x \geq 0} \tilde{F}(x)$$

also is a point of descent for the function  $F(x)$  since

$$F(x_{k+1}) \leq \tilde{F}(x_{k+1}) \leq \tilde{F}(x_k) = F(x_k).$$

Since the constant  $L$  is not known a priori, an inner loop is needed. Algorithm 3 presents an iterative way to carry out this scheme. As in the previous algorithm this also requires only the gradient information and can therefore be applied to two different strategies: to the whole space  $(U, V)$  (Algorithm FFO) and to  $U$  and  $V$  separately in an alternating fashion (Algorithm CFO).

A main difference with the previous algorithm is its stopping criterion for the inner loop. This algorithm requires also a parameter  $\beta$  for which the practical choice is 2.

### Alternative least squares methods

The first algorithm proposed for solving the nonnegative matrix factorization was the alternative least squares method [25]. It is known that, fixing either  $U$  or  $V$ , the problem becomes a least squares problem with nonnegativity constraint.

Since the least squares problems in Algorithm 4 can be perfectly decoupled into smaller problems corresponding to the columns or rows of  $A$ , we can directly apply methods for the Nonnegative Least Square problem to each of the small problem. Methods that can be applied are [19], [6], etc.



**Algorithm 3** (FO)

---

```

1: Initialize  $x^0$ ,  $L_0$  and  $k = 0$ 
2: repeat
3:    $y = [x^k - \frac{1}{L_k} \nabla F(x^k)]_+$ 
4:   while  $F(y) - F(x^k) > \langle \nabla F(x^k), y - x^k \rangle + \frac{L_k}{2} \|y - x^k\|_2^2$  do
5:      $L_k = L_k / \beta$ 
6:      $Y = [x^k - \frac{1}{L_k} \nabla F(x^k)]_+$ 
7:   end while
8:    $x^{k+1} = y$ 
9:    $L_{k+1} = L_k \cdot \beta$ 
10:   $k = k + 1$ 
11: until Stopping condition

```

---

**Algorithm 4** Alternative Least Square (ALS)

---

```

1: Initialize  $U$  and  $V$ 
2: repeat
3:   Solve:  $\min_{V \geq 0} \frac{1}{2} \|A - UV^T\|_F^2$ 
4:   Solve:  $\min_{U \geq 0} \frac{1}{2} \|A^T - VU^T\|_F^2$ 
5: until Stopping condition

```

---

**Implementation**

The most time-consuming job is the test for the sufficient decrease, which is also the stopping condition for the inner loop. As mentioned at the beginning of the section, the above methods can be carried out using two different strategies: full space search or coordinate search. In some cases, it is required to evaluate repeatedly the function  $F(U, V)$ . We mention here how to do this efficiently with the coordinate search.

**Full space search:** The exact evaluation of  $F(x) = F(U, V) = \|A - UV^T\|_F^2$  need  $O(mnr)$  operations. When there is a correction  $y = (U + \Delta U, V + \Delta V)$ , we have to calculate  $F(y)$  which also requires  $O(mnr)$  operations. Hence, it requires  $O(tmnr)$  operations to determine a stepsize in  $t$  iterations of the inner loop.

**Coordinate search:** when  $V$  is fixed, the Euclidean distance is a quadratic function on  $U$ :

$$\begin{aligned}
F(U) = \|A - UV^T\|_F^2 &= \langle A, A \rangle - 2 \langle UV^T, A \rangle + \langle UV^T, UV^T \rangle \\
&= \|A\|_F^2 - 2 \langle U, AV \rangle + \langle U, U(V^T V) \rangle.
\end{aligned}$$

The most expensive step is the computation of  $AV$ , which requires  $O(mnr)$  operations. But when  $V$  is fixed,  $AV$  can be calculated once

at the beginning of the inner loop. The remaining computations are  $\langle U, AV \rangle$  and  $\langle U, U(V^T V) \rangle$ , which requires  $O(nr)$  and  $O(nr^2 + nr)$  operations. Therefore, it requires  $O(tnr^2)$  operations to determine a stepsize in  $t$  iterations of the inner loop which is much less than  $O(tmnr)$  operations. This is due to the assumption  $r \ll n$ . Similarly, when  $U$  fixed,  $O(tmr^2)$  operations are needed to determine a stepsize.

If we consider an iteration is a sweep, i.e. once all the variables are updated, the following table summarizes the complexity of each sweep of the described algorithms:

Algorithm	Complexity per iteration
Mult	$O(mnr)$
FLine	$O(tmnr)$
CLine	$O(t_1nr^2 + t_2mr^2)$
FFO	$O(tmnr)$
CFO	$O(t_1nr^2 + t_2mr^2)$
ALS	$O(2^r mnr)^*$
IALS	$O(mnr)$

where  $t$ ,  $t_1$  and  $t_2$  are the number of iterations of inner loops, which can not be bounded in general. For algorithm *ALS*, the complexity is reported for the case where the active set method [19] is used. Although  $O(2^r mnr)$  is a very high *theoretical* upper bound that count all the possible subsets of  $r$  variables of each subproblem, in practice, the active set method needs much less iterations to converge. One might as well use more efficient convex optimization tools to solve the subproblems instead of the active set method.

## Scaling and Stopping criterion

For descent methods, several stopping conditions are used in the literature. We now discuss some problems when implementing these conditions for NMF.

The very first condition is the decrease of the objective function. The algorithm should stop when it fails to make the objective function decrease with a certain amount :

$$F(U^{k+1}, V^{k+1}) - F(U^k, V^k) < \epsilon \quad \text{or} \quad \frac{F(U^{k+1}, V^{k+1}) - F(U^k, V^k)}{F(U^k, V^k)} < \epsilon.$$

This is not a good choice for all cases since the algorithm may stop at a point very far from a stationary point. Time and iteration bounds can also be imposed for very slowly converging algorithms. But here again this may not be good for the optimality conditions. A better choice is probably the norm of the projected gradient as suggested in [23]. For

the NMF problem it is defined as follows :

$$[\nabla_X^P]_{ij} = \begin{cases} [\nabla_X]_{ij} & \text{if } X_{ij} > 0 \\ \min(0, [\nabla_X]_{ij}) & \text{if } X_{ij} = 0 \end{cases}$$

where  $X$  stands for  $U$  or  $V$ . The proposed condition then becomes

$$\left\| \begin{pmatrix} \nabla_{U^k}^P \\ \nabla_{V^k}^P \end{pmatrix} \right\|_F \leq \epsilon \left\| \begin{pmatrix} \nabla_{U^1} \\ \nabla_{V^1} \end{pmatrix} \right\|_F. \quad (1.18)$$

We should also take into account the scaling invariance between  $U$  and  $V$ . Putting  $\bar{U} = \gamma U$  and  $\bar{V} = \frac{1}{\gamma} V$  does not change the approximation  $UV^T$  but the above projected gradient norm is affected:

$$\begin{aligned} \left\| \begin{pmatrix} \nabla_{\bar{U}}^P \\ \nabla_{\bar{V}}^P \end{pmatrix} \right\|_F^2 &= \|\nabla_{\bar{U}}^P\|_F^2 + \|\nabla_{\bar{V}}^P\|_F^2 = \frac{1}{\gamma^2} \|\nabla_U^P\|_F^2 + \gamma^2 \|\nabla_V^P\|_F^2 \\ &\neq \left\| \begin{pmatrix} \nabla_U^P \\ \nabla_V^P \end{pmatrix} \right\|_F^2. \end{aligned} \quad (1.19)$$

Two approximate factorizations  $UV^T = \bar{U}\bar{V}^T$  resulting in the same approximation should be considered equivalent in terms of precision. One could choose  $\gamma^2 := \|\nabla_U^P\|_F / \|\nabla_V^P\|_F$ , which minimizes (1.19) and forces  $\|\nabla_{\bar{U}}^P\|_F = \|\nabla_{\bar{V}}^P\|_F$ , but this may not be a good choice when only one of the gradients  $\|\nabla_U^P\|_F$  and  $\|\nabla_V^P\|_F$  is nearly zero.

In fact, the gradient  $\begin{pmatrix} \nabla_U \\ \nabla_V \end{pmatrix}$  is scale dependent in the NMF problem and any stopping criterion that uses gradient information is affected by this scaling. To limit that effect, we suggest the following scaling after each iteration:

$$\tilde{U}_k \leftarrow U_k D_k \quad \tilde{V}_k \leftarrow V_k D_k^{-1}$$

where  $D_k$  is a positive diagonal matrix:

$$[D_k]_{ii} = \sqrt{\frac{\|V_{:,i}\|_2}{\|U_{:,i}\|_2}}.$$

This ensures that  $\|\tilde{U}_{:,i}\|_F^2 = \|\tilde{V}_{:,i}\|_F^2$  and hopefully reduces also the difference between  $\|\nabla_{\tilde{U}}^P\|_F^2$  and  $\|\nabla_{\tilde{V}}^P\|_F^2$ . Moreover, it may help to avoid

The same scaling should be applied to the initial point as well ( $U_1, V_1$ ) when using (1.18) as the stopping condition.

## 5. Rank-one Residue Iteration

In the previous section, we have seen that it is very appealing to decouple the problem into convex subproblems. But this may “converge” to solutions that are far from the global minimizers of the problem.

In this section, we analyze a different decoupling of the problem based on rank one approximations. This also allows us to formulate a very simple basic subproblem. This scheme has a major advantage over other methods : the subproblems can be optimally solved in closed form. Therefore it can be proved to have a strong convergence results through its *damped* version and it can be extended to more general types of factorizations such as for nonnegative tensors and to some practical constraints such as sparsity and smoothness. Moreover, the experiments in Section 6 suggest that this method outperforms the other ones in most cases. During the completion of the revised version of this report, we were informed that an independent report [8] had also proposed this decoupling without any convergence investigation and extentions.

## New partition of variables

Let the  $u_i$ 's and  $v_i$ 's be respectively the columns of  $U$  and  $V$ . Then the NMF problem can be rewritten as follows :

**Problem 2** (Nonnegative Matrix Factorization). *Given a  $m \times n$  non-negative matrix  $A$ , solve*

$$\min_{u_i \geq 0, v_i \geq 0} \frac{1}{2} \|A - \sum_{i=1}^r u_i v_i^T\|_F^2.$$

Let us fix all the variables, except for a single vector  $v_t$  and consider the following least squares problem:

$$\min_{v \geq 0} \frac{1}{2} \|R_t - u_t v^T\|_F^2, \quad (1.20)$$

where  $R_t = A - \sum_{i \neq t} u_i v_i^T$ . We have:

$$\|R_t - u_t v^T\|_F^2 = \text{trace} [(R_t - u_t v^T)^T (R_t - u_t v^T)] \quad (1.21)$$

$$= \|R_t\|_F^2 - 2v^T R_t^T u_t + \|u_t\|_2^2 \|v\|_2^2. \quad (1.22)$$

From this formulation, one now derives the following lemma.

**Lemma 5.1.** *If  $[R_t^T u_t]_+ \neq 0$ , then  $v_* := \frac{[R_t^T u_t]_+}{\|u_t\|_2^2}$  is the unique global minimizer of (1.20) and the function value equals  $\|R_t\|_F^2 - \frac{\|[R_t^T u_t]_+\|_2^2}{\|u_t\|_2^2}$ .*

*Proof.* Let us permute the elements of the vectors  $x := R_t^T u_t$  and  $v$  such that

$$Px = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad Pv = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}, \quad \text{with } x_1 \geq 0, \quad x_2 < 0$$

and  $P$  is the permutation matrix. Then

$$\|R_t - u_t v^T\|_F^2 = \|R_t\|_F^2 - 2v_1^T x_1 - 2v_2^T x_2 + \|u_t\|_2^2 (v_1^T v_1 + v_2^T v_2).$$

Since  $x_2 < 0$  and  $v_2 \geq 0$ , it is obvious that  $\|R_t - u_t v^T\|_F^2$  can only be minimal if  $v_2 = 0$ . Our assumption implies that  $x_1$  is nonempty and  $x_1 > 0$ . Moreover  $[R_t^T u_t]_+ \neq 0$  and  $u_t \geq 0$  imply  $\|u_t\|_2^2 > 0$ , one can then find the optimal  $v_1$  by minimizing the remaining quadratic function

$$\|R_t\|_F^2 - 2v_1^T x_1 + \|u_t\|_2^2 v_1^T v_1$$

which yields the solution  $v_1 = \frac{x_1}{\|u_t\|_2^2}$ . Putting the two components together yields the result

$$v_* = \frac{[R_t^T u_t]_+}{\|u_t\|_2^2} \quad \text{and} \quad \|R_t - u_t v_*^T\|_F^2 = \|R_t\|_F^2 - \frac{\|[R_t^T u_t]_+\|_2^2}{\|u_t\|_2^2}.$$

□

---

**Algorithm 5** (RRI)

---

```

1: Initialize  $u_i$ 's,  $v_i$ 's, for  $i = 1$  to  $r$ 
2: repeat
3:   for  $t = 1$  to  $r$  do
4:      $R_t = A - \sum_{i \neq t} u_i v_i^T$ 
5:
6:     if  $[R_t^T u_t]_+ \neq 0$  then
7:        $v_t \leftarrow \frac{[R_t^T u_t]_+}{\|u_t\|_2^2}$ 
8:     else
9:        $v_t = 0$ 
10:    end if
11:
12:    if  $[R_t v_t]_+ \neq 0$  then
13:       $u_t \leftarrow \frac{[R_t v_t]_+}{\|v_t\|_2^2}$ 
14:    else
15:       $u_t = 0$ 
16:    end if
17:  end for
18: until Stopping condition

```

---

**Remark 1:** The above lemma has of course a dual form, where one fixes  $v_t$  but solves for the optimal  $u$  to minimize  $\|R_t - uv_t^T\|_F^2$ . This

would yield the updating rules

$$v_t \leftarrow \frac{[R_t^T u_t]_+}{\|u_t\|_2^2} \quad \text{and} \quad u_t \leftarrow \frac{[R_t v_t]_+}{\|v_t\|_2^2} \quad (1.23)$$

which can be used to recursively update approximations  $\sum_{i=1}^r u_i v_i^T$  by modifying each rank-one matrix  $u_t v_t^T$  in a cyclic manner. This problem is different from the NMF, since the error matrices  $R_t = A - \sum_{i \neq t} u_i v_i^T$  are no longer nonnegative. We will therefore call this method the *Rank-one Residue Iteration* (RRI), i.e. Algorithm 5. The same algorithm was independently reported as Hierarchical Alternating Least Squares (HALS) [8].

**Remark 2:** In case where  $[R_t^T u_t]_+ = 0$ , we have a trivial solution for  $v = 0$  that is not covered by Lemma 5.1. In addition, if  $u_t = 0$ , this solution is no longer unique. In fact,  $v$  can be arbitrarily taken to construct a rank-deficient approximation. The effect of this on the convergence of the algorithm will be discussed further in the next section.

**Remark 3:** Notice that the optimality of Lemma 5.1 implies that  $\|A - UV^T\|$  can not increase. And since  $A \geq 0$  fixed,  $UV^T \geq 0$  must be bounded. Therefore, its component  $u_i v_i^T$  ( $i=1 \dots r$ ) must be bounded as well. One can moreover scale the vector pairs  $(u_i, v_i)$  at each stage as explained in Section 4.0 without affecting the local optimality of Lemma 5.1. It then follows that the rank one products  $u_i v_i^T$  and their scaled vectors remain bounded.

## Convergence

In the previous section, we have established the partial updates for each of the variable  $u_i$  or  $v_i$ . And for a NMF problem where the reduced rank is  $r$ , we have in total  $2r$  vector variables (the  $u_i$ 's and  $v_i$ 's). The described algorithm can be also considered as a projected gradient method since the update (1.23) can be rewritten as:

$$\begin{aligned} u_t &\leftarrow \frac{[R_t v_t]_+}{\|v_t\|_2^2} = \frac{[(A - \sum_{i \neq t} u_i v_i^T) v_t]_+}{\|v_t\|_2^2} = \frac{[(A - \sum_i u_i v_i^T + u_t v_t^T) v_t]_+}{\|v_t\|_2^2} \\ &= \frac{[(A - \sum_i u_i v_i^T) v_t + u_t v_t^T v_t]_+}{\|v_t\|_2^2} = \left[ u_t - \frac{1}{\|v_t\|_2^2} \nabla_{u_t} \right]_+. \end{aligned}$$

Similarly, the update for  $v_i$  can be rewritten as

$$v_t \leftarrow \left[ v_t - \frac{1}{\|u_t\|_2^2} \nabla_{v_t} \right]_+.$$

Therefore, the new method follows the projected gradient scheme described in the previous section. But it produces the optimal solution

in closed form. For each update of a column  $v_t$  (or  $u_t$ ), the proposed algorithm requires just a matrix-vector multiplication  $R_t^T u_t$  (or  $R_t v_t$ ), wherein the residue matrix  $R_t = A - \sum_{i \neq t} u_i v_i^T$  does not have to be calculated explicitly. Indeed, by calculating  $R_t^T u_t$  (or  $R_t v_t$ ) from  $A^T u_t$  (or  $A v_t$ ) and  $\sum_{i \neq t} v_i (u_i^T u_t)$  (or  $\sum_{i \neq t} u_i (v_i^T v_t)$ ), the complexity is reduced from  $O(mnr + mn)$  to only  $O(mn + (m + n)(r - 1))$  which is majored by  $O(mn)$ . This implies that the complexity of each sweep through the  $2r$  variables  $u_t$ 's and  $v_t$ 's requires only  $O(mnr)$  operations, which is equivalent to a sweep of the multiplicative rules and to an inner loop of any gradient methods. This is very low since the evaluation of the whole gradient requires already the same complexity.

Because at each step of the  $2r$  basic steps of Algorithm 5, we compute an optimal rank-one nonnegative correction to the corresponding error matrix  $R_t$  the Frobenius norm of the error can not increase. This is a reassuring property but it does not imply convergence of the algorithm.

Each vector  $u_t$  or  $v_t$  lies in a convex set  $\mathbb{U}_t \subset \mathbb{R}_+^m$  or  $\mathbb{V}_t \subset \mathbb{R}_+^n$ . Moreover, because of the possibility to include scaling we can set an upper bound for  $\|U\|$  and  $\|V\|$ , in such a way that all the  $\mathbb{U}_t$  and  $\mathbb{V}_t$  sets can be considered as closed convex. Then, we can use the following Theorem 5.1, to prove a stronger convergence result for Algorithm 5.

**Theorem 5.1.** *Every limit point generated by Algorithm 5 is a stationary point.*

*Proof.* We notice that, if  $u_t = 0$  and  $v_t = 0$  at some stages of Algorithm 5, they will remain zero and no longer take part in all subsequent iterations. We can divide the execution of Algorithm 5 into two phases.

During the first phase, some of the pairs  $(u_t, v_t)$  become zero. Because there are only a finite number ( $2r$ ) of such vectors, the number of iterations in this phase is also finite. At the end of this phase, we can rearrange and partition the matrices  $U$  and  $V$  such that

$$U = (U_+ \ 0) \text{ and } V = (V_+ \ 0),$$

where  $U_+$  and  $V_+$  do not have any zero column. We temporarily remove zero columns out of the approximation.

During the second phase, no column of  $U_+$  and  $V_+$  becomes zero, which guarantees the updates for the columns of  $U_+$  and  $V_+$  are unique and optimal. Moreover,  $\frac{1}{2} \|A - \sum_{i=1}^r u_i v_i^T\|_F^2$  is continuously differentiable over the set  $\mathbb{U}_1 \times \dots \times \mathbb{U}_r \times \mathbb{V}_1 \times \dots \times \mathbb{V}_r$ , and the  $\mathbb{U}_i$ 's and  $\mathbb{V}_i$ 's are closed convex. A direct application of Proposition 2.7.1 in [4] proves that every stationary point  $(U_+^*, V_+^*)$  is a stationary point. It is then easy to prove that if there are zero columns removed at the end

of the first phase, adding them back yields another stationary point:  $U^* = (U_+^* \ 0)$  and  $V^* = (V_+^* \ 0)$  of the required dimension. However, in this case, the rank of the approximation will then be lower than the requested dimension  $r$ .  $\square$

In Algorithm 5, variables are updated in this order:  $u_1, v_1, u_2, v_2, \dots$ . We can alternate the variables in a different order as well, for example  $u_1, u_2, \dots, u_r, v_1, v_2, \dots, v_r, \dots$ . Whenever this is carried out in a cyclic fashion, the Theorem 5.1 still holds and this does not increase the complexity of each iteration of the algorithm.

As pointed above, stationary points given by Algorithm 5 may contain useless zero components. To improve this, one could replace  $u_t v_t^T (\equiv 0)$  by any nonnegative rank-one approximation that reduces the norm of the error matrix. For example, the substitution

$$u_t = e_{i^*} \quad v_t = [R_t^T u_t]_+, \quad (1.24)$$

where  $i^* = \operatorname{argmax}_i \|[R_t^T e_i]_+\|_2^2$ , reduces the error norm by  $\|[R_t^T e_i]_+\|_2^2 > 0$  unless  $R_t \leq 0$ . These substitutions can be done as soon as  $u_t$  and  $v_t$  start to be zero. If we do these substitutions in only a *finite* number of times before the algorithm starts to converge, Theorem 5.1 still holds. In practice, only a few such substitutions in total are usually needed by the algorithm to converge to a stationary point without any zero component. Note that the matrix rank of the approximation might not be  $r$ , even when all  $u_t$ 's and  $v_t$ 's ( $t = 1 \dots r$ ) are nonzero.

A possibly better way to fix the problem due to zero components is to use the following *damped RRI algorithm* in which we introduce new  $2r$  dummy variables  $w_i \in \mathbb{U}_i$  and  $z_i \in \mathbb{V}_i$ , where  $i = 1 \dots r$ . The new problem to solve is:

**Problem 3** (Damped Nonnegative Matrix Factorization).

$$\min_{\substack{u_i \geq 0 \\ w_i \geq 0}} \min_{\substack{v_i \geq 0 \\ z_i \geq 0}} \frac{1}{2} \|A - \sum_{i=1}^r u_i v_i^T\|_F^2 + \frac{\psi}{2} \sum_i \|u_i - w_i\|_2^2 + \frac{\psi}{2} \sum_i \|v_i - z_i\|_2^2,$$

where the damping factor  $\psi$  is a positive constant.

Again, the coordinate descent scheme is applied with the cyclic update order:  $u_1, w_1, v_1, z_1, u_2, w_2, v_2, z_2, \dots$  to result in the following optimal updates for  $u_t, v_t, w_t$  and  $z_t$ :

$$u_t = \frac{[R_t v_t]_+ + \psi w_t}{\|v_t\|_2^2 + \psi}, \quad w_t = u_t, \quad v_t = \frac{[R_t^T u_t]_+ + \psi z_t}{\|u_t\|_2^2 + \psi} \quad \text{and} \quad z_t = v_t \quad (1.25)$$



where  $t = 1 \dots r$ . The updates  $w_t = u_t$  and  $z_t = v_t$  can be integrated in the updates of  $u_t$  and  $v_t$  to yield Algorithm 6. We have the following results:

**Theorem 5.2.** *Every limit point generated by Algorithm 6 is a stationary point of NMF problem 2.*

---

**Algorithm 6** (Damped RRI)

---

- 1: Initialize  $u_i$ 's,  $v_i$ 's, for  $i = 1$  to  $r$
  - 2: **repeat**
  - 3:   **for**  $t = 1$  to  $r$  **do**
  - 4:      $R_t = A - \sum_{i \neq t} u_i v_i^T$
  - 5:      $v_t \leftarrow \frac{[R_t^T u_t + \psi v_t]_+}{\|u_t\|_2^2 + \psi}$
  - 6:      $u_t \leftarrow \frac{[R_t v_t + \psi u_t]_+}{\|v_t\|_2^2 + \psi}$
  - 7:   **end for**
  - 8: **until** Stopping condition
- 

*Proof.* Clearly the cost function in Problem 3 is continuously differentiable over the set  $\mathbb{U}_1 \times \dots \times \mathbb{U}_r \times \mathbb{U}_1 \times \dots \times \mathbb{U}_r \times \mathbb{V}_1 \times \dots \times \mathbb{V}_r \times \mathbb{V}_1 \times \dots \times \mathbb{V}_r$ , and the  $\mathbb{U}_i$ 's and  $\mathbb{V}_i$ 's are closed convex. The uniqueness of the global minimum of the elementary problems and a direct application of Proposition 2.7.1 in [4] prove that every limit point of Algorithm 6 is a stationary point of Problem 3.

Moreover, at a stationary point of Problem 3, we have  $u_t = w_t$  and  $v_t = z_t$ ,  $t = 1 \dots r$ . The cost function in Problem 3 becomes the cost function of the NMF problem 2. This implies that every stationary point of Problem 3 yields a stationary point of the standard NMF problem 2.  $\square$

This *damped* version not only helps to eliminate the problem of zero components in the convergence analysis but may also help to avoid zero columns in the approximation when  $\psi$  is carefully chosen. But it is not an easy task. Small values of  $\psi$  provide an automatic treatment of zeros while not changing much the updates of RRI. Larger values of  $\psi$  might help to prevent the vectors  $u_t$  and  $v_t$  ( $t = 1 \dots r$ ) from becoming zero too soon. But too large values of  $\psi$  limit the updates to only small changes, which will slow down the convergence.

In general, the rank of the approximation can still be lower than the requested dimension. Patches may still be needed when a zero component appears. Therefore, in our experiments, using the *undamped* RRI algorithm 5 with the substitution (1.24) is still the best choice.

## Variants of the RRI method

We now extend the Rank-one Residue Iteration by using a factorization of the type  $XDY^T$  where  $D$  is diagonal and nonnegative and the columns of the nonnegative matrices  $X$  and  $Y$  are normalized. The NMF formulation then becomes

$$\min_{\substack{x_i \in \mathbb{X}_i, y_i \in \mathbb{Y}_i \\ d_i \in \mathbb{R}_+}} \frac{1}{2} \|A - \sum_{i=1}^r d_i x_i y_i^T\|_F^2,$$

where  $\mathbb{X}_i$ 's and  $\mathbb{Y}_i$ 's are sets of normed vectors.

The variants that we present here depend on the choice of  $\mathbb{X}_i$ 's and  $\mathbb{Y}_i$ 's. A generalized Rank-one Residue Iteration method for low-rank approximation is given in Algorithm 7. This algorithm needs to solve a sequence of elementary problems of the type:

$$\max_{s \in \mathbb{S}} y^T s \quad (1.26)$$

where  $y \in \mathbb{R}^n$  and  $\mathbb{S} \subset \mathbb{R}^n$  is a set of normed vectors. We first introduce a permutation vector  $I_y = (i_1 \ i_2 \ \dots \ i_n)$  which reorders the elements of  $y$  in non-increasing order :  $y_{i_k} \geq y_{i_{k+1}}$ ,  $k = 1 \dots (n-1)$ . The function  $p(y)$  returns the number of positive entries of  $y$ .

---

### Algorithm 7 GRRI

---

- 1: Initialize  $x_i$ 's,  $y_i$ 's and  $d_i$ 's, for  $i = 1$  to  $r$
  - 2: **repeat**
  - 3:   **for**  $i = 1$  to  $r$  **do**
  - 4:      $R_i = A - \sum_{j \neq i} d_j x_j y_j^T$
  - 5:      $y_i \leftarrow \operatorname{argmax}_{s \in \mathbb{Y}_i} (x_i^T R_i s)$
  - 6:      $x_i \leftarrow \operatorname{argmax}_{s \in \mathbb{X}_i} (y_i^T R_i^T c)$
  - 7:      $d_i = x_i^T R_i y_i$
  - 8:   **end for**
  - 9: **until** Stopping condition
- 

Let us first point out that for the set of normed nonnegative vectors the solution of problem (1.26) is given by  $s^* = \frac{y_+}{\|y_+\|_2}$ . It then follows that Algorithm 7 is essentially the same as Algorithm 5 since the solutions  $v_i$  and  $u_i$  of each step of Algorithm 7, given by (1.23), correspond exactly to those of problem (1.26) via the relations  $y_i = u_i / \|u_i\|_2$ ,  $y_i = v_i / \|v_i\|_2$  and  $d_i = \|u_i\|_2 \|v_i\|_2$ .

Below we list the sets for which the solution  $s^*$  of (1.26) can be easily computed.

- *Set of normed vectors:*  $s = \frac{y}{\|y\|_2}$ . This is useful when one wants to create factorizations where only one of the factor  $U$  or  $V$  is nonnegative and the other is real matrix.
- *Set of normed nonnegative vectors:*  $s = \frac{y_+}{\|y_+\|_2}$ .
- *Set of normed bounded nonnegative vectors  $\{s\}$ :* where  $0 \leq l_i \leq s_i \leq p_i$ . The optimal solution of (1.26) is given by:

$$s = \max \left( l, \min \left( p, \frac{y_+}{\|y_+\|_2} \right) \right).$$

- *Set of normed binary vectors  $\{s\}$ :* where  $s = \frac{b}{\|b\|}$  and  $b \in \{0, 1\}^n$ . The optimal solution of (1.26) is given by:

$$[s^*]_{i_t} = \begin{cases} \frac{1}{\sqrt{k^*}} & \text{if } t \leq k^* \\ 0 & \text{otherwise} \end{cases} \quad \text{where } k^* = \underset{k}{\operatorname{argmax}} \frac{\sum_{t=1}^k y_{i_t}}{\sqrt{k}}.$$

- *Set of normed sparse nonnegative vectors:* all normed nonnegative vectors having at most  $K$  nonzero entries. The optimal solution for (1.26) is given by norming the following vector  $p^*$

$$[p^*]_{i_t} = \begin{cases} y_{i_t} & \text{if } t \leq \min(p(y), K) \\ 0 & \text{otherwise} \end{cases}$$

- *Set of normed fixed-sparsity nonnegative vectors:* all nonnegative vectors  $s$  a fixed sparsity, where

$$\text{sparsity}(s) = \frac{\sqrt{n} - \|s\|_1 / \|s\|_2}{\sqrt{n} - 1}.$$

The optimal solution for (1.26) is given by using the projection scheme in [17].

One can also imagine other variants, for instance by combining the above ones. Depending on how data need to be approximated, one can create new algorithms provided it is relatively simple to solve problem (1.26). There have been some particular ideas in the literatures such as NMF with sparseness constraint [17], Semidiscrete Matrix Decomposition [18] and Semi-Nonnegative Matrix Factorization [11] for which variants of the above scheme can offer an alternative choice of algorithm.

**Remark:** Only the first three sets are the normed version of a closed convex set, as required for the convergence by Theorem 5.1. Therefore the algorithms might not converge to a stationary point with the other sets. However, the algorithm always guarantees a non-increasing update even in those cases and can therefore be expected to return a *good* approximation.

## Nonnegative Tensor Factorization

If we refer to the problem of finding the nearest nonnegative vector to a given vector  $a$  as the nonnegative approximation in one dimension, the NMF is its generalization in two dimensions and naturally, it can be extended to even higher-order tensor approximation problems. Algorithms described in the previous section use the closed form solution of the one dimensional problem to solve the two-dimensional problem. We now generalize this to higher orders. Since in one dimension such an approximation is easy to construct, we continue to use this approach to build the solutions for higher order problems.

For a low-rank tensor, there are two popular kinds of factored tensors, namely those of Tucker and Kruskal [2]. We only give an algorithm for finding approximations of Kruskal type. It is easy to extend this to tensors of Tucker type, but this is omitted here.

Given a  $d$  dimensional tensor  $T$ , we will derive an algorithm for approximating a nonnegative tensor by a rank- $r$  nonnegative Kruskal tensor  $S \in \mathbb{R}_+^{n_1 \times n_2 \times \dots \times n_d}$  represented as a sum of  $r$  rank-one tensors:

$$S = \sum_{i=1}^r \sigma_i u_{1i} \star u_{2i} \star \dots \star u_{di}$$

where  $\sigma_i \in \mathbb{R}_+$  is a scaling factor,  $u_{ti} \in \mathbb{R}_+^{n_t}$  is a normed vector (i.e.  $\|u_{ti}\|_2 = 1$ ) and  $a \star b$  stands for the outer product between two vectors or tensors  $a$  and  $b$ .

The following update rules are the generalization of the matrix case to the higher order tensor:

$$y = (\dots ((\dots (R_k u_{1k}) \dots u_{(t-1)k}) u_{(t+1)k}) \dots) u_{dk} \quad (1.27)$$

$$\sigma_k = \|[y]_+\|_2, \quad u_{tk} = \frac{[y]_+}{\sigma_k}, \quad (1.28)$$

where  $R_k = T - \sum_{i \neq k} \sigma_i u_{1i} \star u_{2i} \star \dots \star u_{di}$  is the residue tensor calculated without the  $k^{th}$  component of  $S$  and  $R_k u_{ij}$  is the ordinary tensor/vector product in the corresponding dimension.

We can then produce an algorithm which updates in a cyclic fashion all vectors  $u_{ji}$ . This is in fact a direct extension to Algorithm 5, one can carry out the same discussion about the convergence here to guarantee that each limit point of this algorithm is a stationary point for the non-negative tensor factorization problem and to improve the approximation quality.

Again, as we have seen in the previous section, we can extend the procedure to take into account different constraints on the vectors  $u_{ij}$  such as discreteness, sparseness, etc.

The approach proposed here is again different from that in [9] where a similar cascade procedure for multilayer nonnegative matrix factorization is used to compute a 3D tensor approximation. Clearly, the approximation error will be higher than our proposed method, since the cost function is minimized by taking into account all the dimensions.

## Regularizations

The regularizations are common methods to cope with the ill-posedness of inverse problems. Having known some additional information about the solution, one may want to impose a priori some constraints to algorithms, such as: smoothness, sparsity, discreteness, etc. To add such regularizations in to the RRI algorithms, it is possible to modify the NMF cost function by adding some regularizing terms. We will list here the update for  $u_i$ 's and  $v_i$ 's when some simple regularizations are added to the original cost function. The proof of these updates are straightforward and hence omitted.

- One-Norm  $\|\cdot\|_1$  regularization: the one-norm of the vector variable can be added as a heuristic for finding a sparse solution. This is an alternative to the fixed-sparsity variant presented above. The regularized cost function with respect to the variable  $v_t$  will be

$$\frac{1}{2}\|R_t - u_t v^T\|_F^2 + \beta\|v\|_1, \quad \beta > 0$$

where the optimal update is given by

$$v_t^* = \frac{[R_t^T u_t - \beta \mathbf{1}_{n \times 1}]_+}{\|u_t\|_2^2}.$$

The constant  $\beta > 0$  can be varied to control the trade-off between the approximation error  $\frac{1}{2}\|R_t - u_t v^T\|_F^2$  and  $\|v\|_1$ . From this update, one can see that this works by zeroing out elements of  $R_t^T u_t$  which are smaller than  $\beta$ , hence reducing the number of nonzero elements of  $v_t^*$ .

- Smoothness regularization  $\|v - B\hat{v}_t\|_F^2$ : where  $\hat{v}_t$  is the current value of  $v_t$  and the matrix  $B$  helps to calculate the average of the neighboring elements at each element of  $v$ . When  $v$  is a 1D smooth

function,  $B$  can be the following  $n \times n$  matrix:

$$B = \begin{pmatrix} 0 & 1 & \dots & \dots & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & \dots & 0 & 1 & 0 \end{pmatrix}. \quad (1.29)$$

This matrix can be defined in a different way to take the true topology of  $v$  into account, for instance  $v = \text{vec}(F)$  where  $F$  is a matrix. The regularized cost function with respect to the variable  $v_t$  will be

$$\frac{1}{2}\|R_t - u_t v^T\|_F^2 + \frac{\delta}{2}\|v - B\hat{v}_t\|_F^2, \quad \delta > 0$$

where the optimal update is given by

$$v_t^* = \frac{[R_t^T u_t + \delta B\hat{v}_t]_+}{\|u_t\|_2^2 + \delta}.$$

The constant  $\delta \geq 0$  can be varied to control the trade-off between the approximation error  $\frac{1}{2}\|R_t - u_t v^T\|_F^2$  and the smoothness of  $v_t$  at the fixed point. From the update, one can see that this works by searching for the optimal update  $v_t^*$  with some preference for the neighborhood of  $B\hat{v}_t$ , i.e., a smoothed vector of the current value  $\hat{v}_t$ .

The two above regularizations can be added independently to each of the columns of  $U$  and/or  $V$ . The trade-off factor  $\beta$  (or  $\delta$ ) can be different for each column. A combination of different regularizations on a column (for instance  $v_t$ ) can also be used to solve the multi-criterion problem

$$\frac{1}{2}\|R_t - u_t v^T\|_F^2 + \frac{\gamma}{2}\|v\|_2^2 + \frac{\delta}{2}\|v - B\hat{v}_t\|_F^2, \quad \beta, \gamma, \delta > 0$$

where the optimal update is given by

$$v_t^* = \frac{[R_t^T u_t - \beta \mathbf{1}_{n \times 1} + \delta B\hat{v}_t]_+}{\|u_t\|_2^2 + \delta}.$$

The one-norm regularizations as well as the two-norm regularization can be found in [1] and [3]. A major difference with that method is that the norm constraints is added to the rows rather than on the columns of  $V$  or  $U$  as done here. However, for the two versions of the one-norm

regularization, the effects are somehow similar. While the two-norm regularization on the columns of  $U$  and  $V$  are simply scaling effects, which yield nothing in the RRI algorithm. We therefore only test the smoothness regularization at the end of the chapter with some numerical generated data.

For more extensions and variants, see [16].

## 6. Experiments

Here we present several experiments to compare the different descent algorithms presented in this paper. For all the algorithms, the scaling scheme proposed in section 4.0 was applied.

### Random matrices

We generated 100 random nonnegative matrices of different sizes. We used seven different algorithms to approximate each matrix:

- the multiplicative rule (**Mult**),
- alternative least squares using Matlab function *lsqnonneg* (**ALS**),
- a full space search using line search and Armijo criterion (**FLine**),
- a coordinate search alternating on  $U$  and  $V$ , and using line search and Armijo criterion (**CLine**),
- a full space search using first-order approximation (**FFO**),
- a coordinate search alternating on  $U$  and  $V$ , and using first-order approximation (**CFO**),
- an iterative rank-one residue approximation (**RRI**).

For each matrix, the same starting point is used for every algorithm. We create a starting point by randomly generating two matrices  $U$  and  $V$  and then rescaling them to yield a first approximation of the original matrix  $A$  as proposed in Section 4.0:

$$U = UD\sqrt{\alpha}, \quad V = VD^{-1}\sqrt{\alpha},$$

where

$$\alpha := \frac{\langle A, UV^T \rangle}{\langle UV^T, UV^T \rangle} \quad \text{and} \quad D_{ij} = \begin{cases} \sqrt{\frac{\|V_{:i}\|_2}{\|U_{:i}\|_2}} & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}.$$

From (1.15), we see that when approaching a KKT stationary point of the problem, the above scaling factor  $\alpha \rightarrow 1$ . This implies that every KKT stationary point of this problem is scale-invariant.

$\epsilon$	Mult	ALS	FLine	CLine	FFO	CFO	RRI
(m=30, n=20, r=2)							
$10^{-2}$	0.02(96)	0.40	0.04	0.02	0.02	0.01	0.01
$10^{-3}$	0.08(74)	1.36	0.12	0.09	0.05	0.04	0.03
$10^{-4}$	0.17(71)	2.81	0.24	0.17	0.11	0.08	0.05
$10^{-5}$	0.36(64)	4.10	0.31	0.25	0.15	0.11	0.07
$10^{-6}$	0.31(76)	4.74	0.40	0.29	0.19	0.15	0.09
(m=100, n=50, r=5)							
$10^{-2}$	45*(0)	3.48	0.10	0.09	0.09	0.04	0.02
$10^{-3}$	45*(0)	24.30(96)	0.59	0.63	0.78	0.25	0.15
$10^{-4}$	45*(0)	45*(0)	2.74	2.18	3.34	0.86	0.45
$10^{-5}$	45*(0)	45*(0)	5.93	4.06	6.71	1.58	0.89
$10^{-6}$	45*(0)	45*(0)	7.23	4.75	8.98	1.93	1.30
(m=100, n=50, r=10)							
$10^{-2}$	45*(0)	11.61	0.28	0.27	0.18	0.11	0.05
$10^{-3}$	45*(0)	41.89(5)	1.90	2.11	1.50	0.74	0.35
$10^{-4}$	45*(0)	45*(0)	7.20	5.57	5.08	2.29	1.13
$10^{-5}$	45*(0)	45*(0)	12.90	9.69	10.30	4.01	1.71
$10^{-6}$	45*(0)	45*(0)	14.62(99)	11.68(99)	13.19	5.26	2.11
(m=100, n=50, r=15)							
$10^{-2}$	45*(0)	25.98	0.66	0.59	0.40	0.20	0.09
$10^{-3}$	45*(0)	45*(0)	3.90	4.58	3.18	1.57	0.61
$10^{-4}$	45*(0)	45*(0)	16.55(98)	13.61(99)	9.74	6.12	1.87
$10^{-5}$	45*(0)	45*(0)	21.72(97)	17.31(92)	16.59(98)	7.08	2.39
$10^{-6}$	45*(0)	45*(0)	25.88(89)	19.76(98)	19.20(98)	10.34	3.66
(m=100, n=100, r=20)							
$10^{-2}$	45*(0)	42.51(4)	1.16	0.80	0.89	0.55	0.17
$10^{-3}$	45*(0)	45*(0)	9.19	8.58	10.51	5.45	1.41
$10^{-4}$	45*(0)	45*(0)	28.59(86)	20.63(94)	29.89(69)	12.59	4.02
$10^{-5}$	45*(0)	45*(0)	32.89(42)	27.94(68)	34.59(34)	18.83(90)	6.59
$10^{-6}$	45*(0)	45*(0)	37.14(20)	30.75(60)	36.48(8)	22.80(87)	8.71
(m=200, n=100, r=30)							
$10^{-2}$	45*(0)	45*(0)	2.56	2.20	2.68	1.31	0.44
$10^{-3}$	45*(0)	45*(0)	22.60(99)	25.03(98)	29.67(90)	12.94	4.12
$10^{-4}$	45*(0)	45*(0)	36.49(2)	39.13(13)	45*(0)	33.33(45)	14.03
$10^{-5}$	45*(0)	45*(0)	45*(0)	39.84(2)	45*(0)	37.60(6)	21.96(92)
$10^{-6}$	45*(0)	45*(0)	45*(0)	45*(0)	45*(0)	45*(0)	25.61(87)

Table 1.1. Comparison of average successful running time of algorithms over 100 random matrices. Time limit is 45 seconds. 0.02(96) means that a result is returned with the required precision  $\epsilon$  within 45 seconds for 96 (of 100) matrices of which the average running time is 0.02 seconds. 45\*(0): failed in all 100 matrices.



The algorithms are all stopped when the projected gradient norm is lower than  $\epsilon$  times the gradient norm at the starting point or when it takes more than 45 seconds. The relative precisions  $\epsilon$  are chosen equal to  $10^{-2}$ ,  $10^{-3}$ ,  $10^{-4}$ ,  $10^{-5}$ ,  $10^{-6}$ . No limit was imposed on the number of iterations.

For alternative gradient algorithms **CLine** and **CFO**, we use different precisions  $\epsilon_U$  and  $\epsilon_V$  for each of the inner iteration for  $U$  and for  $V$  as suggested in [23] where  $\epsilon_U$  and  $\epsilon_V$  are initialized by  $10^{-3}$ . And when the inner loop for  $U$  or  $V$  needs no iteration to reach the precision  $\epsilon_U$  or  $\epsilon_V$ , one more digit of precision will be added into  $\epsilon_U$  or  $\epsilon_V$  (i.e.  $\epsilon_U = \epsilon_U/10$  or  $\epsilon_V = \epsilon_V/10$ ).

Table 1.1 shows that for all sizes and ranks, Algorithm **RRI** is the fastest to reach the required precision. Even though it is widely used in practice, algorithm **Mult** fails to provide solutions to the NMF problem within the allocated time. A further investigation shows that the algorithm gets easily trapped in boundary points where some  $U_{ij}$  and/or  $V_{ij}$  is zero while  $\nabla_{U_{ij}}$  and/or  $\nabla_{V_{ij}}$  is negative, hence violating one of the KKT conditions (1.11). The multiplicative rules then fail to move and do not return to a local minimizer. A slightly modified version of this algorithm was given in [21], but it needs to wait to get sufficiently close to such points before attempting an escape, and is therefore also not efficient. The **ALS** algorithm can return a stationary point, but it takes too long.

We select five methods: **FLine**, **CLine**, **FFO**, **CFO** and **RRI** for a more detailed comparison. For each matrix  $A$ , we run these algorithms with 100 different starting points. Figure 1.1, 1.2, 1.3 and 1.4 show the results with some different settings. One can see that, when the approximated errors are almost the same between the algorithms, **RRI** is the best overall in terms of running times. It is probably because the RRI algorithm chooses only one vector  $u_t$  or  $v_t$  to optimize at once. This allows the algorithm to move *optimally* down on partial direction rather than just a *small step* on a more global direction. Furthermore, the computational load for an update is very small, only one matrix-vector multiplication is needed. All these factors make the running time of the RRI algorithm very attractive.

## Image data

The following experiments use the Cambridge ORL face database as the input data. The database contains 400 images of 40 persons (10 images per person). The size of each image is  $112 \times 92$  with 256 gray levels per pixel representing a front view of the face of a person. The images

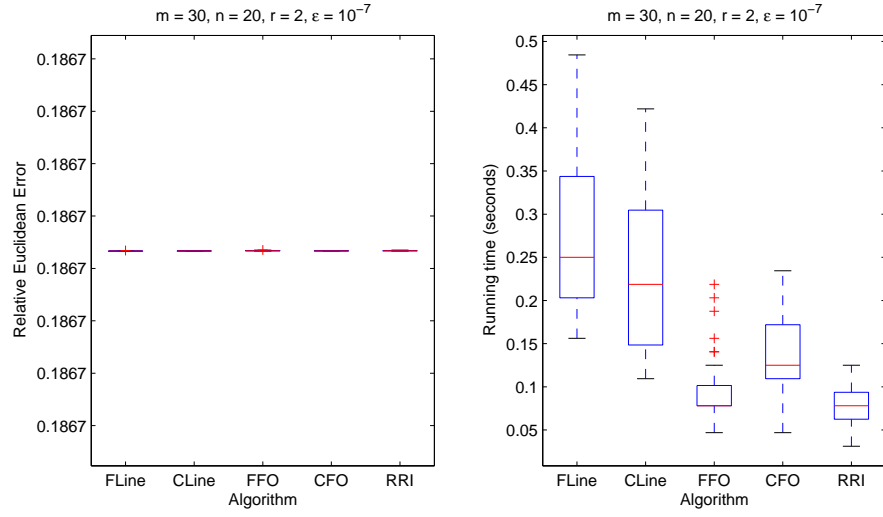


Figure 1.1. Comparison of selected algorithms for  $\epsilon = 10^{-7}$

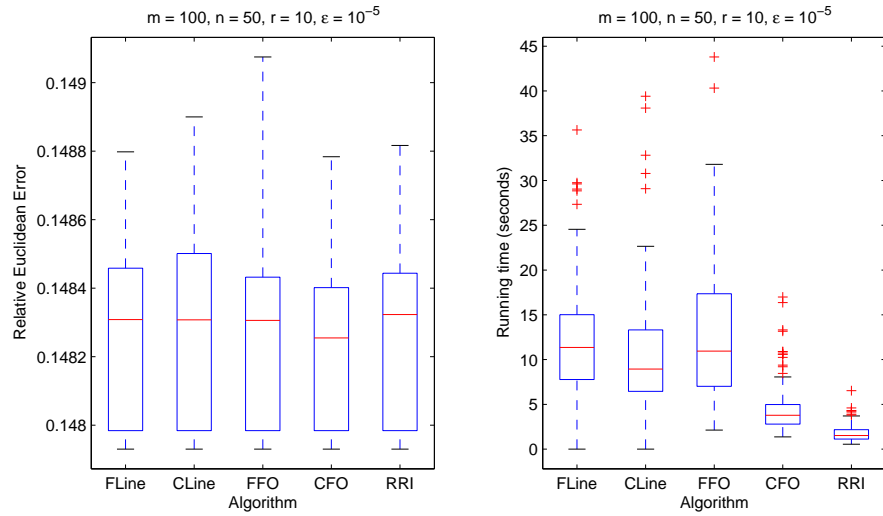
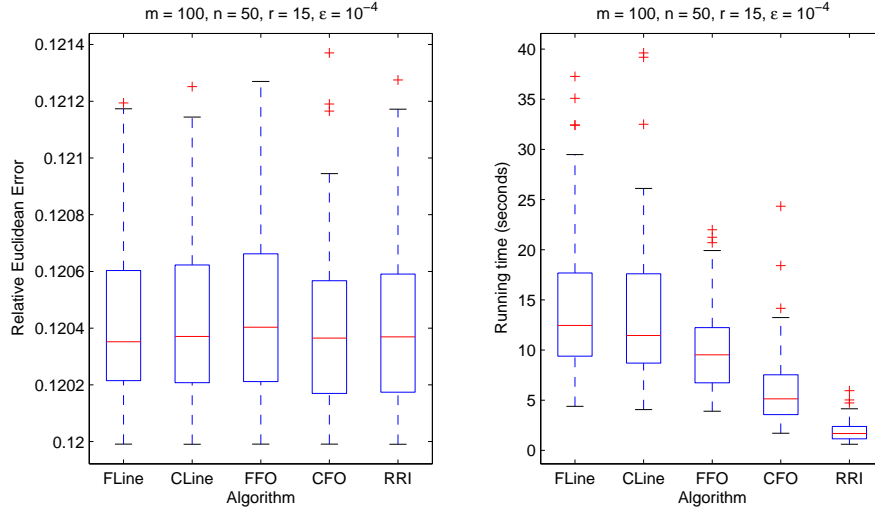
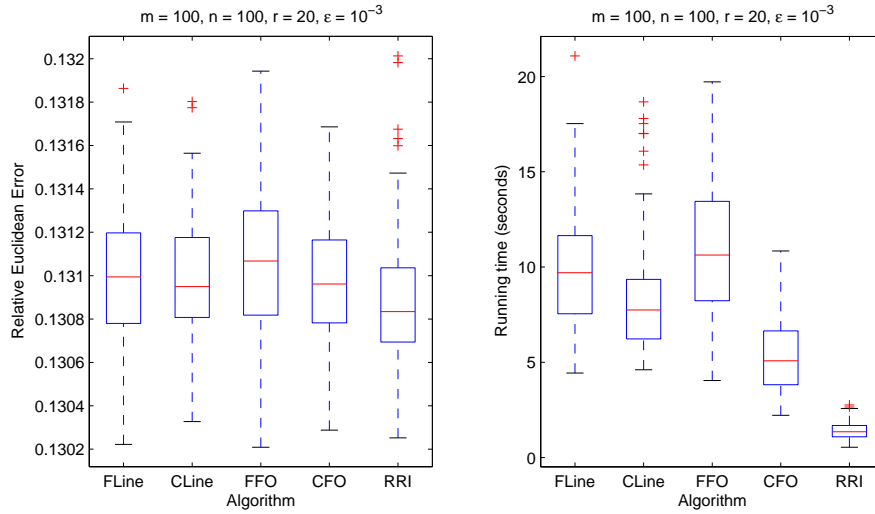


Figure 1.2. Comparison of selected algorithms for  $\epsilon = 10^{-5}$

Figure 1.3. Comparison of selected algorithms for  $\epsilon = 10^{-4}$ Figure 1.4. Comparison of selected algorithms for  $\epsilon = 10^{-3}$

are then transformed into 400 “face vectors” in  $\mathbb{R}^{10304}$  ( $112 \times 92 = 10304$ ) to form the data matrix  $A$  of size  $10304 \times 400$ . We used three weight matrices of the same size of  $A$  (ie.  $10304 \times 400$ ). Since it was used in [20], this data has become the standard benchmark for NMF algorithms.

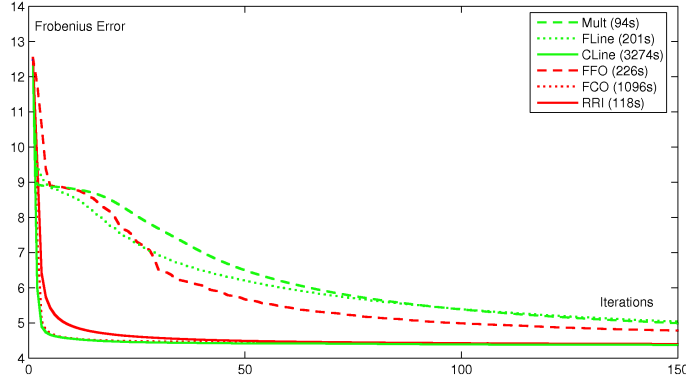


Figure 1.5. NMF: Error vs. Iterations

In the first experiment, we run six NMF algorithms described above on this data for the reduced rank of 49. The original matrix  $A$  is constituted by transforming each image into one of its column. Figure 1.5 shows for the six algorithms the evolution of the error versus the number of iterations. Because the minimization process is different in each algorithm, we will say that one iteration corresponds to all elements of both  $U$  and  $V$  being updated. Figure 1.6 shows the evolution of the error versus time. Since the work of one iteration varies from one algorithm to another, it is crucial to plot the error versus time to get a fair comparison between the different algorithms. In the two figures, we can see that the RRI algorithm behaves very well on this dataset. And since its computation load of each iteration is small and constant (without inner loop), this algorithm converges faster than the others.

In the second experiment, we construct a third-order nonnegative tensor approximation. We first build a tensor by stacking all 400 images to have a  $112 \times 92 \times 400$  nonnegative tensor. Using the proposed algorithm, a  $rank-142$  nonnegative tensor is calculated to approximate this tensor. Figure 1.7 shows the result for six images chosen randomly from the 400 images. Their approximations given by the  $rank-142$  nonnegative tensor are much better than that given by the  $rank-8$  nonnegative matrix, even though they require similar storage space:  $8 * (112 * 92 + 400) = 85632$  and  $142 * (112 + 92 + 400) = 85768$ . The  $rank-8$  truncated SVD approximation (i.e.  $[A_8]_+$ ) is also included for reference.

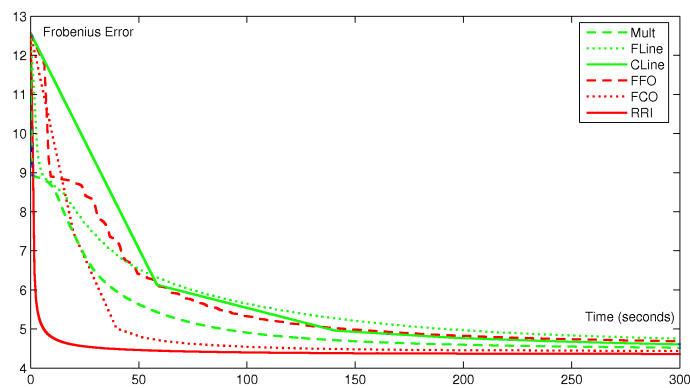


Figure 1.6. NMF: Error vs. Time

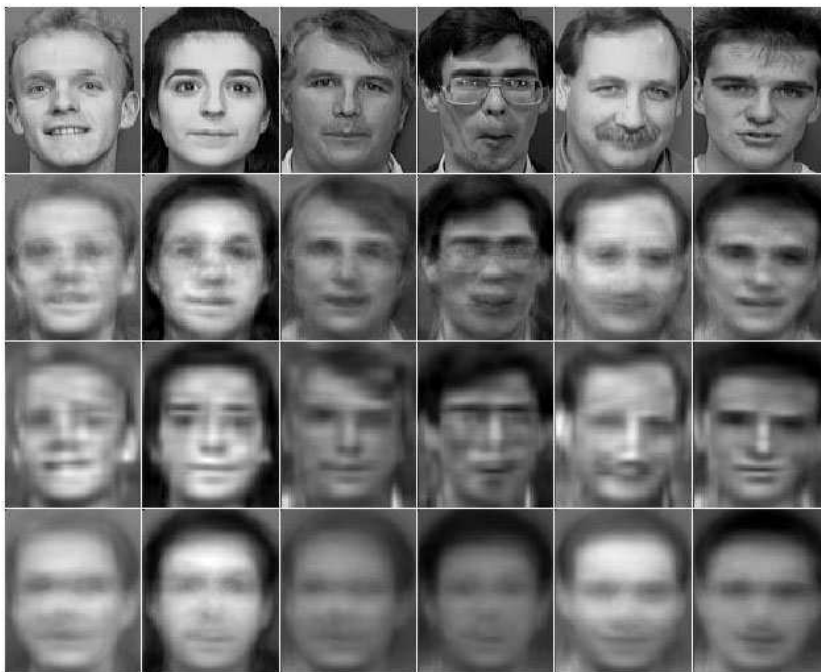


Figure 1.7. Tensor Factorization vs. Matrix Factorization on facial data. Six randomly chosen images from 400 of ORL dataset. From top to bottom: original images, their  $\text{rank} - 8$  truncated SVD approximation, their  $\text{rank} - 142$  nonnegative tensor approximation (150 RRI iterations) and their  $\text{rank} - 8$  nonnegative matrix approximation (150 RRI iterations).

In the third experiment, we apply the variants of RRI algorithm mentioned in Section 5.0 to the face databases. The following settings are compared:

- **Original**: original faces from the databases.
- **49NMF**: standard factorization (nonnegative vectors),  $r = 49$ .
- **100Binary**: columns of  $U$  are limited to the scaled binary vectors,  $r = 100$ .
- **49Sparse10**: columns of  $U$  are sparse. Not more than 10% of the elements of each column of  $A$  are positive.  $r = 49$ .
- **49Sparse20**: columns of  $U$  are sparse. Not more than 20% of the elements of each column of  $A$  are positive.  $r = 49$ .
- **49HSparse60**: columns of  $U$  are sparse. The Hoyer sparsity of each column of  $U$  are 0.6.  $r = 49$ .
- **49HSparse70**: columns of  $U$  are sparse. The Hoyer sparsity of each column of  $U$  are 0.7.  $r = 49$ .
- **49HBSparse60**: columns of  $U$  are sparse. The Hoyer sparsity of each column of  $U$  are 0.6. Columns of  $V$  are scaled binary.  $r = 49$ .
- **49HBSparse70**: columns of  $U$  are sparse. The Hoyer sparsity of each column of  $U$  are 0.7. Columns of  $V$  are scaled binary.  $r = 49$ .

For each setting, we use RRI algorithm to compute the corresponding factorization. Some randomly selected faces are reconstructed by these settings as shown in Figure 1.8. For each setting, RRI algorithm produces a different set of bases to approximate the original faces. When the columns of  $V$  are constrained to scaled binary vectors (**100Binary**), the factorization can be rewritten as  $UV^T = \hat{U}B^T$ , where  $B$  is a binary matrix. This implies that each image is reconstructed by just the presence or absence of 100 bases shown in Figure 1.9.

Figure 1.10 and 1.11 show nonnegative bases obtained by imposing some sparsity on the columns of  $V$ . The sparsity can be easily controlled by the percentages of positive elements or by the Hoyer sparsity measure.

Figure 1.12 combines the sparsity of the bases (columns of  $U$ ) and the binary representation of  $V$ . The sparsity is measured by the Hoyer measure as in Figure 1.11. Only with the absence or presence of these 49 *features*, faces are approximated as showed in the last two rows of Figure 1.8.

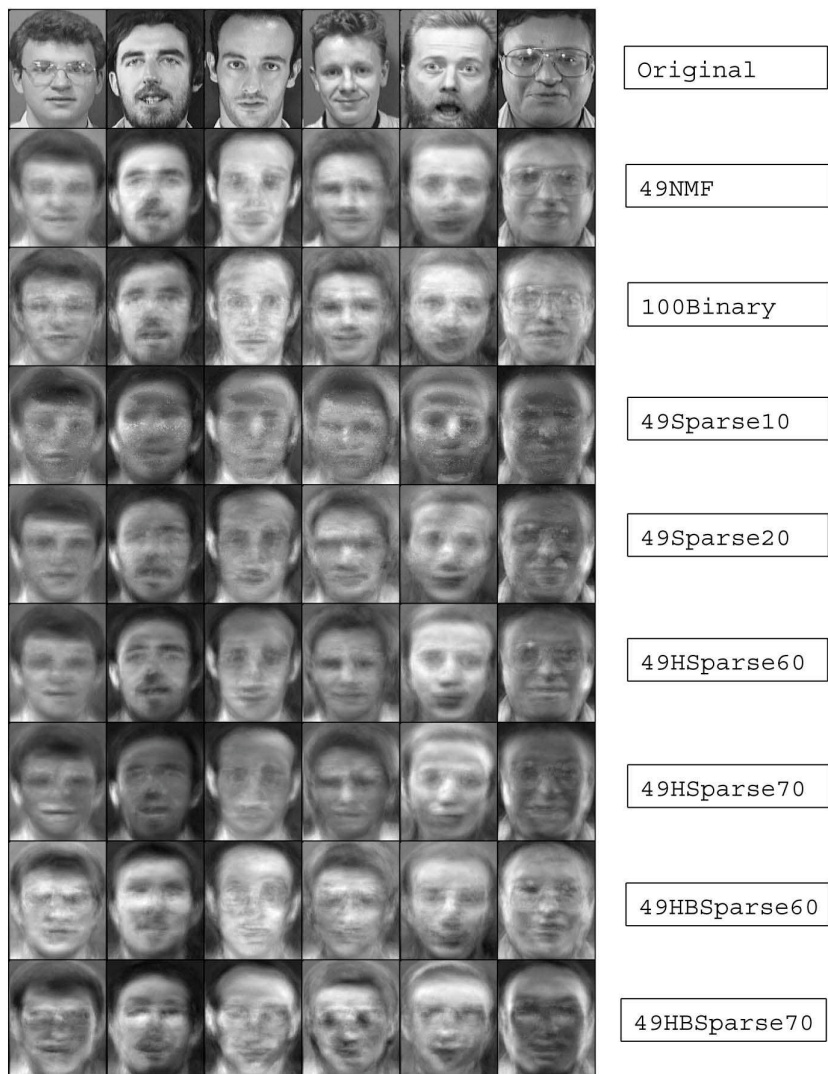


Figure 1.8. Nonnegative matrix factorization with several sparse settings



Figure 1.9. Bases from **100Binary** setting



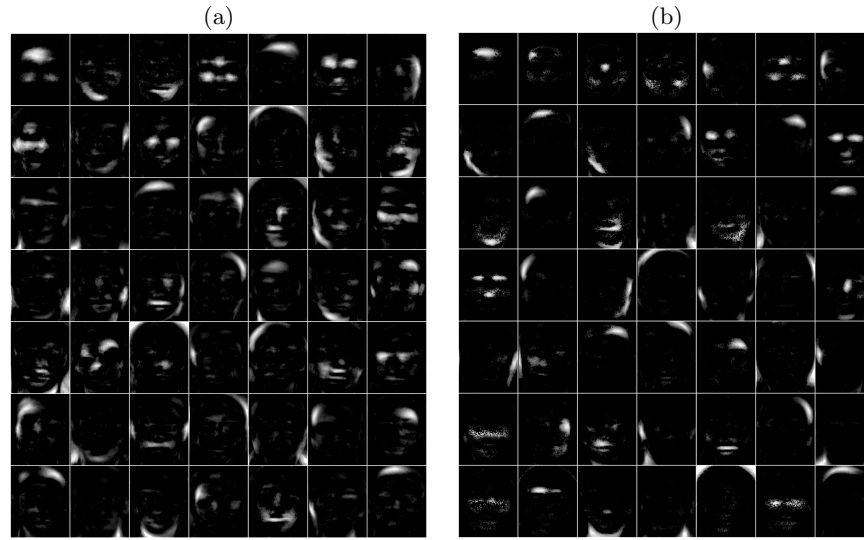


Figure 1.10. Sparse bases **49Sparse20** and **49Sparse10**. Maximal percentage of positive elements is 20% (a) and 10% (b)

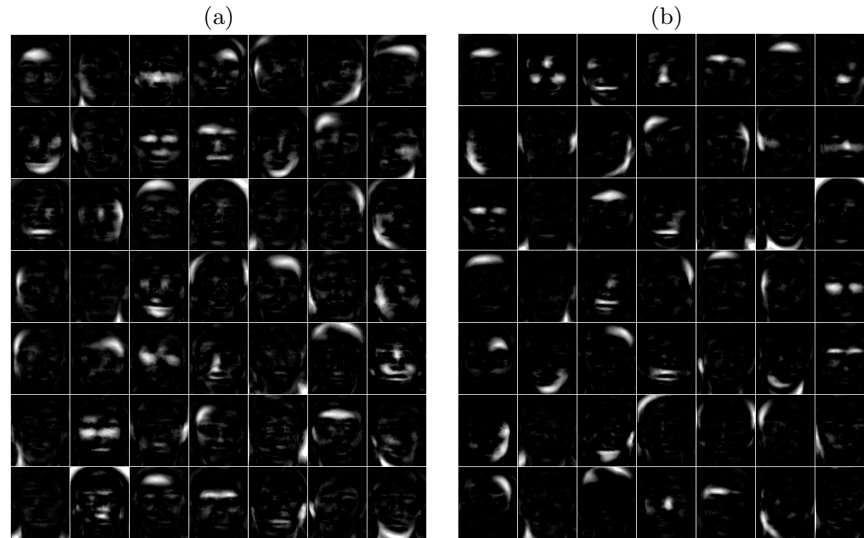


Figure 1.11. Hoyer sparse bases **49HSparse60** and **49HSparse70**. Sparsity of bases is 0.6 (a) and 0.7 (b)

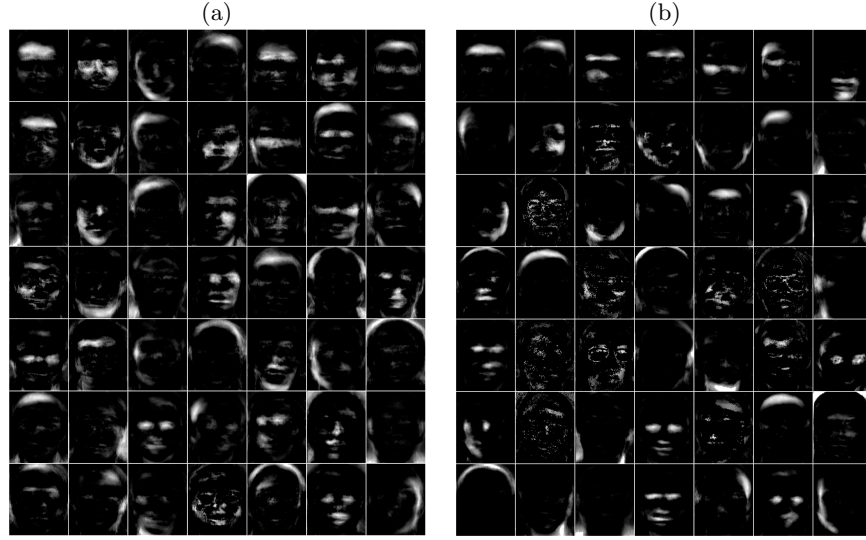


Figure 1.12. Hoyer sparse bases **49HBSparse60** and **49HBSparse70**. Sparsity of bases is 0.6 (a) and 0.7 (b).  $V$  is binary matrix.

The above examples show how to use the variants of the RRI algorithm to control the sparsity of the bases. One can see that the sparser the bases are, the less storage is needed to store the approximation. Moreover, this provides a part-based decomposition using *local* features of the faces.

### Smooth approximation

We carry out this experiment to test the new smoothness constraint introduced in the previous section:

$$\frac{1}{2}\|R_i - u_i v^T\|_F^2 + \frac{\delta}{2}\|v - B\hat{v}_i\|_F^2, \quad \delta > 0$$

where  $B$  is defined in (1.29).

We generate the data using four smooth nonnegative functions  $f_1, f_2, f_3$  et  $f_4$ , described in Figure 1.13, where each function is represented as a nonnegative vector of size 200.

We then generate a matrix  $A$  containing 100 mixture of these functions as follows

$$A = \max(FE^T + N, 0)$$

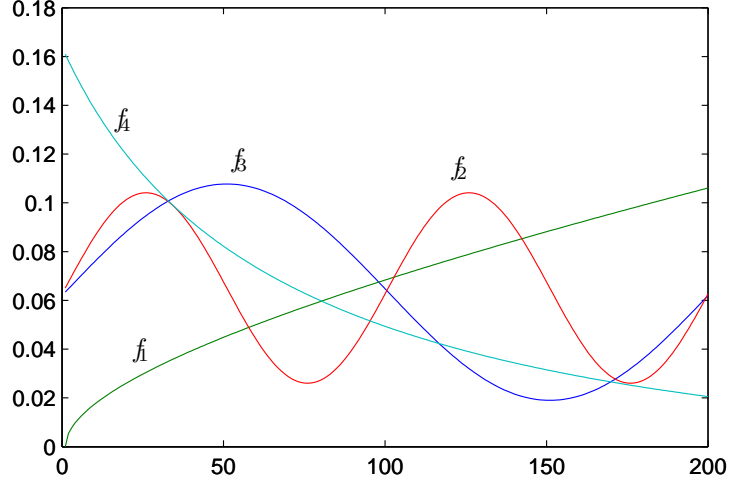


Figure 1.13. Smooth functions

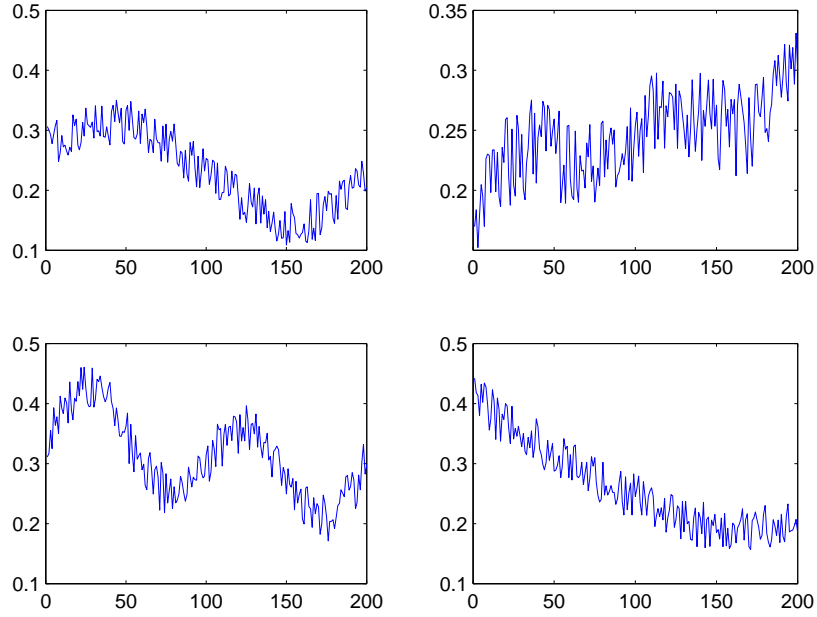


Figure 1.14. Randomly selected generated data

where  $F = [f_1 \ f_2 \ f_3 \ f_4]$ ,  $E$  is a random nonnegative matrix and  $N$  is normally distributed random noise with  $\|N\|_F = 0.2\|FE^T\|_F$ . Four randomly selected columns of  $A$  are plotted in Figure 1.14.

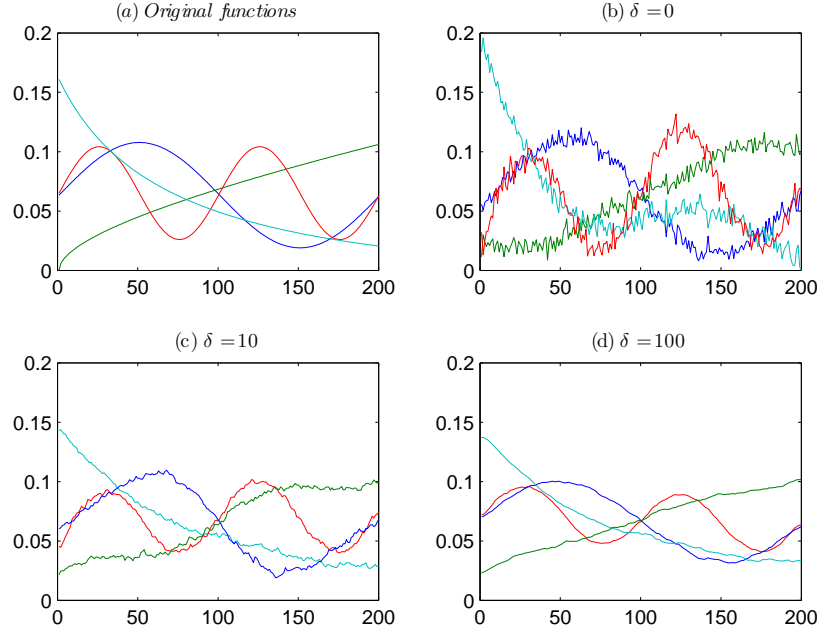


Figure 1.15. Original functions vs. reconstructed functions

We run the regularized RRI algorithm to force the smoothness of columns of  $U$ . We apply, for each run, the same value of  $\delta$  for all the columns of  $U$ :  $\delta = 0, 10, 100$ . The results obtained through these runs are presented in Figure 1.15. We see that, without regularization, i.e.  $\delta = 0$ , the noise is present in the approximation, which produces nonsmooth solutions. When increasing the regularizing terms, i.e.  $\delta = 10, 100$ , the reconstructed functions become smoother and the shape of the original functions are well preserved.

This smoothing technique can be used for applications like that in [27], where smooth spectral reflectance data from space objects is unmixed. The multiplicative rules are modified by adding the two-norm regularizations on the factor  $U$  and  $V$  to enforce the smoothness. This is a different approach, therefore, a comparison should be carried out.

We have described a new method for nonnegative matrix factorization that has a good and fast convergence. Moreover, it is also very flexible to create variants and to add some constraints as well. The numerical experiments show that this method and its derived variants behave very well with different types of data. This gives enough motivations to extend to other types of data and applications in the future. In the last

two chapters of this thesis, it is applied to weighted cost functions and to symmetric factorizations.

## 7. Conclusion

This paper focuses on the descent methods for Nonnegative Matrix Factorization, which are characterized by nonincreasing updates at each iteration.

We present also the Rank-one Residue Iteration algorithm for computing an approximate Nonnegative Matrix Factorization. It uses recursively nonnegative rank one approximations of a residual matrix that is not necessarily nonnegative. This algorithm requires no parameter tuning, has nice properties and typically converges quite fast. It also has many potential extensions. During the revision of this report, we were informed that essentially the same algorithm was published in an independent contribution [8] and also mentioned later in an independent personal communication [12].

## Acknowledgments

This paper presents research results of the Concerted Research Action(ARC) "Large Graphs and Networks" of the French Community of Belgium and the Belgian Network DYSCO (Dynamical Systems, Control, and Optimization), funded by the Interuniversity Attraction Poles Programme, initiated by the Belgian State, Science Policy Office. The scientific responsibility rests with its authors. Ngoc-Diep Ho is a FRIA fellow.

## References

- [1] R. Albright, J. Cox, D. Duling, A.N. Langville, and C.D. Meyer. Algorithms, initializations, and convergence for the nonnegative matrix factorization. *Preprint*, 2006.
- [2] B.W. Bader and T.G. Kolda. Efficient MATLAB computations with sparse and factored tensors. Technical Report SAND2006-7592, Sandia National Laboratories, Albuquerque, NM and Livermore, CA, Dec. 2006.
- [3] M.W. Berry, M. Browne, A.N. Langville, V.P. Pauca, and R.J. Plemmons. Algorithms and applications for approximate non-negative matrix factorization. *Computational Statistics and Data Analysis*, 52(1):155–173, 2007.
- [4] D.P. Bertsekas. *Nonlinear programming*. Athena Scientific Belmont, Mass, 1999.

- [5] M. Biggs, A. Ghodsi and S. Vavasis. Nonnegative matrix factorization via rank-one downdate. *University of Waterloo, Preprint*, 2007.
- [6] R. Bro and S. De Jong. A fast non-negativity constrained least squares algorithm. *Journal of Chemometrics*, 11(5):393–401, 1997.
- [7] M. Catral, L. Han, M. Neumann, and R.J. Plemmons. On reduced rank nonnegative matrix factorization for symmetric nonnegative matrices. *Linear Algebra and Its Applications*, 393:107–126, 2004.
- [8] A. Cichocki, R. Zdunek, and S. Amari. Hierarchical ALS Algorithms for Nonnegative Matrix and 3D Tensor Factorization. In *Proceedings of Independent Component Analysis, ICA 2007, London, UK, September 9-12, 2007, Lecture Notes in Computer Science*, Springer, 4666:169–176, 2007.
- [9] A. Cichocki, R. Zdunek, and S. Amari. Nonnegative matrix and tensor factorization *IEEE on Signal Processing Magazine*, 25:142–145, 2008.
- [10] A. Cichocki, R. Zdunek. NMFLAB for Signal Processing, available at <http://www.bsp.brain.riken.jp/ICALAB/nmflab.html>.
- [11] C. Ding, T. Li, and M.I. Jordan. Convex and Semi-Nonnegative Matrix Factorizations. Technical report, LBNL Tech Report 60428, 2006.
- [12] N. Gillis and Francois Glineur. Nonnegative Matrix Factorization and Underapproximation. Preprint, 2007.
- [13] G. Golub and C.F. Van Loan. *Matrix computations.3rd ed.* Baltimore, The Johns Hopkins Univ. Press. xxvii, 694 p. , 1996.
- [14] N.J. Higham, M.J.C. Gover and S. Barnett. Matrix Nearness Problems and Applications. *Applications of Matrix Theory*, 1–27, 1989.
- [15] N.-D. Ho, P. Van Dooren, and V.D. Blondel. Descent algorithms for Nonnegative Matrix Factorization. *Technical Report 2007-57, Cesame. University catholique de Louvain*. Belgium. 2007.
- [16] N.-D. Ho. Nonnegative Matrix Factorization - Algorithms and Applications. *PhD Thesis. University catholique de Louvain*. Belgium. 2008.
- [17] P.O. Hoyer. Non-negative matrix factorization with sparseness constraints. *Journal of Machine Learning Research*, 5:1457–1469, 2004.
- [18] T.G. Kolda and D.P. O’Leary. A semidiscrete matrix decomposition for latent semantic indexing information retrieval. *ACM Transactions on Information Systems (TOIS)*, 16(4):322–346, 1998.

- [19] C.L. Lawson and R.J. Hanson. *Solving least squares problems*. Prentice-Hall Englewood Cliffs, NJ, 1974.
- [20] D.D. Lee and H.S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401:788–791, 1999.
- [21] C.-J. Lin. On the convergence of multiplicative update algorithms for non-negative matrix factorization. *IEEE Transactions on Neural Networks*, 2007. To appear.
- [22] C.-J. Lin. Projected gradient methods for non-negative matrix factorization. *Neural Computation*, 2007. To appear.
- [23] C.-J.-Lin. Projected gradient methods for non-negative matrix factorization. *Technical Report Information and Support Services Technical Report ISSTECH-95-013, Department of Computer Science, National Taiwan University, 2005*
- [24] M. Merritt and Y. Zhang. Interior-Point Gradient Method for Large-Scale Totally Nonnegative Least Squares Problems. *Journal of Optimization Theory and Applications*, 126(1):191–202, 2005.
- [25] P. Paatero and U. Tapper. Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5(1):111–126, 1994.
- [26] P. Paatero. A weighted non-negative least squares algorithm for three-way 'parafac' factor analysis. *Chemometrics and Intelligent Laboratory Systems*, 38(2):223–242, 1997.
- [27] V. P. Pauca, J. Piper and R. J. Plemmons. Nonnegative matrix factorization for spectral data analysis *Linear Algebra and its Applications*, 416(1):29–47, 2006.